

Содержание

Предисловие	9
Глава 1. Элегантный NumPy: фундамент научного программирования на Python	32
Введение в данные: что такое экспрессия гена?	34
N-мерные массивы NumPy	38
Зачем использовать массивы ndarray вместо списков Python?	39
Векторизация	41
Транслирование	41
Исследование набора данных экспрессии генов	43
Чтение данных при помощи библиотеки pandas	43
Нормализация	46
Нормализация между образцами	46
Нормализация между генами	52
Нормализация по образцам и генам: RPKM	54
Подведение итогов	61
Глава 2. Квантильная нормализация с NumPy и SciPy	62
Получение данных	64
Разница в распределении экспрессии генов между индивидуумами	65
Бикластеризация количественных данных	68
Визуализация кластеров	70
Предсказание выживаемости	72
Дальнейшая работа: использование кластеров пациентов TCGA	77
Дальнейшая работа: воспроизведение кластеров TCGA	77
Глава 3. Создание сетей из областей изображений при помощи ndimage	78
Изображения – это просто массивы NumPy	79
Задача: добавление сеточного наложения	84
Фильтры в обработке сигналов	84
Фильтрация изображений (двумерные фильтры)	90
Универсальные фильтры: произвольные функции от соседних значений	92
Задача: игра «Жизнь» Конуэя	93
Задача: магнитуа градиента Собела	94
Графы и библиотека NetworkX	94
Задача: подбор кривой при помощи SciPy	98

Графы смежности областей.....	98
Элегантный пакет ndimage: как строить графы из областей изображений ...	102
Собираем все вместе: сегментация по среднему цвету	105
Глава 4. Частота и быстрое преобразование Фурье	107
Введение в частоту	107
Иллюстрация: спектрограмма пения птиц	110
История	115
Реализация	115
Выбор длины ДПФ	116
Дополнительные понятия ДПФ	118
Частоты и их упорядочивание	118
Оконное преобразование.....	124
Практическое применение: анализ радарных данных.....	128
Свойства сигнала в частотной области	133
Оконное преобразование на практике	136
Радарные изображения.....	138
Дополнительные применения БПФ	142
Дополнительные материалы для чтения.....	143
Задача: свертывание изображения	143
Глава 5. Таблицы сопряженности на основе разреженных координатных матриц	144
Таблицы сопряженности.....	146
Задача: вычислительная сложность матриц ошибок.....	147
Задача: альтернативный алгоритм вычисления матрицы ошибок.....	147
Задача: мультиклассовая матрица ошибок	148
Форматы данных модуля scipy.sparse	148
Формат COO	148
Задача: представление в формате COO	149
Формат сжатой разреженной строки	150
Применения разреженных матриц: преобразования изображений	152
Задача: поворот изображения	156
Назад к таблицам сопряженности	157
Задача: сокращение объема потребляемой оперативной памяти	158
Таблицы сопряженности в сегментации изображений	159
Теория информации вкратце	160
Задача: вычисление условной энтропии	163
Теория информации применительно к сегментации: изменчивость информации.....	163
Конвертирование программного кода массивов NumPy под использование разреженных матриц	166

Применение изменчивости информации	167
Дальнейшая работа: сегментация на практике.....	173
Глава 6. Линейная алгебра в SciPy	174
Основы линейной алгебры	174
Лапласова матрица графа	175
Задача: матрица поворота	176
Лапласовы матрицы с данными о мозге.....	181
Задача: изображение аффинного подобия.....	186
Задача: линейная алгебра с разреженными матрицами	186
PageRank: линейная алгебра для репутации и важности	187
Задача: обработка висячих узлов	192
Задача: эквивалентность разных методов получения собственного вектора	192
Заключительные замечания	192
Глава 7. Оптимизация функций в SciPy	193
Оптимизация в SciPy: <code>scipy.optimize</code>	195
Пример: вычисление оптимального сдвига изображения.....	195
Регистрация изображения при помощи <code>optimize</code>	201
Предотвращение локальных минимумов на основе алгоритма <code>basin hopping</code>	204
Задача: модификация функции <code>align</code>	205
«Что лучше?»: выбор правильной целевой функции.....	205
Глава 8. Большие данные с Toolz в маленьком ноутбуке	212
Потоковая передача при помощи <code>yield</code>	214
Введение в потоковую библиотеку Toolz	217
Подсчет k-мер и исправление ошибок.....	219
Каррирование: изюминка потоковой обработки.....	223
Возвращаясь к подсчету k-мер	226
Задача: анализ главных компонент потоковых данных.....	227
Марковская модель на основе полного генома	228
Задача: онлайн-овая распаковка архива	231
Эпилог	233
Что дальше?.....	233
Списки рассылок	233
GitHub	234
Конференции	235
За пределами SciPy	235
Содействие этой книге	236

До следующей встречи.....	237
Приложение. Решения задач.....	238
Решение: добавление сеточного наложения	238
Решение: игра «Жизнь» Конуэя»	239
Решение: магнитуа градиента Собела	240
Решение: подбор кривой при помощи SciPy	241
Решение: свертывание изображения.....	243
Решение: вычислительная сложность матриц ошибок	243
Решение: альтернативный алгоритм вычисления матрицы ошибок.....	243
Решение: вычисление матрицы ошибок	244
Решение: представление в формате COO	244
Решение: поворот изображения.....	245
Решение: сокращение объема потребляемой оперативной памяти	246
Решение: вычисление условной энтропии.....	247
Решение: матрица поворота.....	247
Решение: изображение аффинного подобия.....	248
Решение: линейная алгебра с разреженными матрицами.....	249
Решение: обработка висячих узлов.....	252
Решение: методы проверки	253
Решение: модификация функции align	253
Решение: анализ главных компонент потоковых данных при помощи библиотеки scikit-learn	255
Решение: добавление шага в начало конвейера	257
Предметный указатель	259

Предисловие

В отличие от стереотипного подвечного стиля, оно было – если использовать технический термин – элегантным, как компьютерный алгоритм, который всего несколькими строками исходного кода достигает впечатляющего результата.

– Грэм Симсион, «Проект “Рози”»

Добро пожаловать в книгу «*Элегантный SciPy*». Мы собираемся провести довольно много времени, сосредоточившись на той части заголовка книги, которая относится к «SciPy», поэтому давайте воспользуемся моментом, чтобы поразмышлять над словом «элегантный». Существует масса руководств, учебных пособий и веб-сайтов документации, которые дают всестороннее описание пакета SciPy. Книга «*Элегантный SciPy*» идет дальше. Она представляет собой нечто большее, чем просто обучение приемам написания по-настоящему рабочего программного кода. Мы вдохновим вас на написание программного кода, который будет по-настоящему потрясающим!

В романе «Проект “Рози”» (между прочим, уморительная книга; когда закончите читать «*Элегантный SciPy*», обязательно прочтите в Википедии¹ материал, описывающий предысторию ее написания) Грэм Симсион переиначивает общепринятые нормы, связанные со словом «элегантный». Большинство людей использует это слово для описания визуальной простоты, стиля и изящества, например, первого мобильного iPhone. Вместо этого герой Грэма Симсиона, Дон Тиллман, дает *определение* элегантности, используя термин «компьютерный алгоритм». Мы надеемся, что после прочтения этой книги вы получите ясное понимание того, что он имеет в виду, и что впредь, занимаясь чтением или написанием куска элегантного программного кода, вы будете ощущать умиротворение в лучах его красоты и изящества. (Возьмите на заметку: авторы могут быть подвержены гиперболе.)

Хороший фрагмент программного кода вызывает ощущение удовлетворенности. Когда вы на него смотрите, его замысел *ясен*, он нередко *краток* (но не настолько краток, чтобы быть туманным), и он *эффективен* при выполнении практической задачи. Для авторов удовольствие от анализа элегантного программного кода лежит в скрытых внутри него уроках и в том, как он вдохновляет нас быть *творческими* в подходах к новым алгоритмическим задачам.

Как ни странно, креативность, помимо всего прочего, может также заставлять нас умничать за счет читателя и писать маловразумительный про-

¹ См. https://en.wikipedia.org/wiki/The_Rosie_Project и https://ru.wikipedia.org/wiki/Проект_«Рози».

граммный код, который очень трудно понять. PEP8¹ (руководство по стилю программирования на Python) и PEP20² (дзэн языка Python) нам напоминают, что «программный код читается намного больше раз, чем пишется» и поэтому «читаемость имеет значение».

Краткость элегантного кода обеспечивается не за счет упаковки кучи вложенных вызовов функций, а за счет абстракции и рационального использования функций. Она требует одной-двух минут на то, чтобы вникнуть, но в конечном счете обеспечивает вам четкий момент понимания. Как только вы начнете разбираться в различных компонентах программного кода, его правильность должна стать очевидной. Этому способствуют ясные имена переменных и функций и тщательно продуманные комментарии, которые программный код *объясняют*, а не просто его *описывают*.

Инженер-программист Дж. Брэдфорд Хиппс (J. Bradford Hipps) недавно в газете «Нью-Йорк таймс»³ заявил: «для того чтобы писать хороший программный код, следует почитать Вирджинию Вульф»:

На практике разработка программного обеспечения имеет гораздо более творческий характер, чем алгоритмический.

Разработчик обращается к своему редактору исходного кода таким же образом, как писатель к чистому листку бумаги. [...] Разработчика и писателя также могут отличать общее здоровое нетерпение по отношению к тому, как все «делалось всегда», и воспроизводящееся стремление нарушать общепринятые правила. Когда модуль закончен или страницы завершены, их качество оценивается в сопоставлении со многими из тех же самых стандартов: элегантностью, краткостью, целостностью; обнаружением симметрий там, где их существование ни разу не было замечено. И даже красотой.

Это именно та позиция, которую мы примем в этой книге.

Теперь, когда мы рассмотрели «элегантную» часть заголовка, давайте возвратимся к «SciPy».

В зависимости от контекста «SciPy» может означать пакет программного обеспечения, экосистему или сообщество разработчиков. Отчасти величие SciPy обусловлено тем, что он имеет превосходную онлайн-документацию⁴ и учебные пособия⁵, поэтому предлагать читателю очередной справочник было бы бессмысленно. Вместо этого в книге «Элегантный SciPy» будет представлен самый лучший программный код, который был создан при помощи SciPy.

¹ См. <https://www.python.org/dev/peps/pep-0008/>.

² См. <https://www.python.org/dev/peps/pep-0020/>.

³ См. <https://www.nytimes.com/2016/05/22/opinion/sunday/to-write-software-read-novels.html>.

⁴ См. <https://docs.scipy.org/>.

⁵ См. <http://www.scipy-lectures.org/>.

Отобранный нами программный код подчеркивает умное, элегантное использование расширенного функционала NumPy, SciPy и других связанных с ними библиотек. Начинающий читатель научится применять эти библиотеки к реальным задачам, используя красивый программный код. При этом в обоснование наших примеров мы используем реальные научные данные.

Как и сам SciPy, мы хотели построить книгу «*Элегантный SciPy*» на основе сообщества разработчиков. Многие приводимые в книге примеры были взяты из рабочего программного кода, обнаруженного в обширной экосистеме научного программирования на Python, и отобраны для демонстрации кратко очерченных выше принципов элегантного программного кода.

Для кого эта книга предназначена?

Цель книги «*Элегантный SciPy*» – побудить вас поднять свои навыки программирования на языке Python на более высокий уровень. Вы изучите пакет SciPy на примере самого лучшего программного кода.

Прежде чем приступить к работе, вы должны иметь представление о языке Python и знать, что такое переменные, функции, циклы, и, возможно, немного разбираться в библиотеке NumPy. Хотя вполне может быть, что вы даже отточили свои навыки программирования на Python на основе такого продвинутого материала, как, например, «Python. К вершинам мастерства» (Fluent Python)¹. Если это к вам не относится, то, прежде чем продолжить чтение этой книги, вам следует начать с каких-нибудь учебных пособий по Python для начинающих, таких как Software Carpentry².

Но, возможно, для вас «стек SciPy» – все равно, что пункт меню сети экспресс-блинных ИНОР, и вы чувствуете себя не в своей тарелке, когда речь идет о его применении на практике. Или, возможно, вы являетесь ученым, который прочитал в Сети несколько учебных пособий по Python и скачал несколько аналитических сценариев из другой лаборатории или у другого сотрудника вашей собственной лаборатории и повозился с ними некоторое время. И, возможно, думаете, что в той или иной степени одиноки в своей попытке научиться программировать SciPy. Вы ошибаетесь.

По ходу изложения мы научим вас использовать Интернет в качестве своего справочника. И мы будем направлять вас к спискам рассылок, хранилищам и конференциям, где вы встретите аналогично мыслящих ученых, которые в своем опыте работы находятся немного дальше, чем вы.

Это такая книга, которую вы прочтете один раз, но будете к ней возвращаться в дальнейшем в поисках вдохновения (и, возможно, чтобы еще раз восхищаться некоторыми элегантными фрагментами программного кода!).

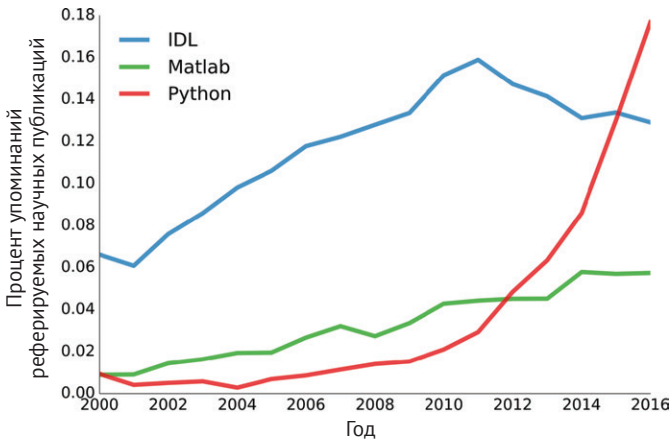
¹ См. <http://shop.oreilly.com/product/0636920032519.do>.

² См. <http://software-carpentry.org/>.

ПОЧЕМУ ИМЕННО SciPy?

Библиотеки NumPy и SciPy составляют ядро научной экосистемы языка Python. В программной библиотеке SciPy реализован набор функций для обработки научных данных из таких областей, как статистика, обработка сигналов, обработка изображений и математическая оптимизация. Библиотека SciPy надстроена поверх библиотеки NumP, которая предназначена для вычислительной обработки числовых массивов. За последние несколько лет вся экосистема приложений и библиотек продемонстрировала существенный рост, опираясь как раз на NumPy и SciPy, с охватом широкого спектра дисциплин, который среди прочих включает астрономию, биологию, метеорологию, климатологию и материаловедение.

И этот рост не проявляет никаких признаков к ослаблению. В 2014 г. Томас Робитэйл и Крис Бомон (Thomas Robitaille и Chris Beaumont) задокументировали¹ рост применения языка Python в астрономии. Вот то, что мы обнаружили, когда обновили² их график во второй половине 2016 г.:



Совершенно очевидно, что в течение многих последующих лет SciPy и связанные с ней библиотеки будут основными в подавляющей части аналитической обработки научных данных.

Еще одним убедительным примером служит тот факт, что организация Software Carpentry, которая обучает ученых вычислительным навыкам, используя для этого чаще всего язык Python, сегодня едва справляется со спросом на свои учебные курсы.

¹ См. https://nbviewer.jupyter.org/github/ChrisBeaumont/adass_proceedings/blob/master/Mining%20acknowledgments%20in%20ADS.ipynb.

² См. <https://gist.github.com/jni/3339985a016572f178d3c2f18e27ec0d>.

ЧТО ТАКОЕ ЭКОСИСТЕМА SciPy?

SciPy (произносится как «Сай Пи») – это экосистема программного обеспечения с открытым исходным кодом на основе Python для математики, науки и техники.

– <http://www.scipy.org>

Экосистема SciPy представляет собой нестрого определенную коллекцию пакетов Python. В книге «Элегантный SciPy» мы встретимся с большинством ее главных компонентов:

- **NumPy** – фундамент научных вычислений на Python. Эта библиотека обеспечивает эффективные числовые массивы и широкую поддержку численных вычислений, включая линейную алгебру, случайные числа и разложение в ряды Фурье. Уникальная особенность NumPy заключена в ее «N-мерных массивах», или массивах ndarray. Эти структуры данных эффективным образом хранят числовые величины и задают решетку в любом количестве размерностей (подробнее об этом чуть позже)¹;
- **SciPy**, как библиотека, является коллекцией эффективных численных алгоритмов для таких областей, как обработка сигналов, интеграция, оптимизация и статистика. Эти алгоритмы обернуты в легкие для использования интерфейсы²;
- **Matplotlib** – мощный пакет для построения графиков в двух измерениях (и элементарных 3D-графиков). Он берет свое название от навевшего его разработку синтаксиса Matlab³;
- **IPython** – интерактивный интерфейс для языка Python, который позволяет оперативно взаимодействовать с вашими данными и тестировать свои идеи⁴;
- блокнот **Jupyter** работает в вашем браузере и позволяет составлять документы с широкими функциональными возможностями, которые объединяют в себе программный код, текст, математические выражения и интерактивные элементы интерфейса⁵. На самом деле при подготовке настоящей книги текст программ был преобразован в блокноты Jupyter и выполнялся там (благодаря этому мы знаем, что все примеры выполняются правильно). Проект Jupyter, начавшийся как расширение IPython,

¹ См. <http://www.numpy.org/>.

² См. <http://www.scipy.org/scipylib/index.html>.

³ См. <http://matplotlib.org/>.

⁴ См. <https://ipython.org/>.

⁵ См.: Перес Ф. «Грамотные вычисления» и вычислительная воспроизводимость: IPython в эпоху журналистики, ориентированной на данные (публикация в блоге). 19 апреля 2013 г. (Fernando Perez. 'Literate computing' and computational reproducibility: IPython in the age of data-driven journalism. <http://blog.fperez.org/2013/04/literate-computing-and-computational.html>).

теперь поддерживает многочисленные языки, включая Cython, Julia, R, Octave, Bash, Perl и Ruby¹;

- **pandas** обеспечивает быстрые столбчатые структуры данных в простом для применения пакете. Он в особенности подходит для работы с помеченными наборами данных, такими как таблицы или реляционные базы данных, и для управления данными временных рядов и скользящими окнами. Кроме того, пакет **pandas** располагает несколькими удобными инструментами, предназначенными для анализа данных, в т. ч. разбором, очисткой и агрегированием данных, а также построением графиков²;
- **scikit-learn** предоставляет унифицированный интерфейс к алгоритмам машинного обучения³;
- **scikit-image** обеспечивает инструменты анализа изображений, которые напрямую интегрируются в остальную часть экосистемы SciPy⁴.

Помимо перечисленных выше библиотек, существует целый ряд других пакетов Python, которые являются частью экосистемы SciPy, и некоторые из них мы также увидим. Хотя в центре внимания настоящей книги будут библиотеки NumPy и SciPy, именно большое количество окружающих их пакетов делает Python движущей силой научных вычислений.

ВЕЛИКИЙ КАТАКЛИЗМ: PYTHON 2 ПРОТИВ PYTHON 3

В ваших путешествиях по Python вы, вероятно, уже слышали пересуды о том, какая версия Python лучше. Возможно, вы задавались вопросом, а разве нельзя просто взять последнюю версию. (Сразу отвечаем: можно.)

В конце 2008 г. разработчики ядра Python выпустили Python 3, который среди прочих улучшений стал основным обновлением языка с более оптимальной обработкой текста (на многочисленных естественных языках) в кодировке Юникод, согласованностью типов и обработкой потоковых данных. Как язвительно заметил Дуглас Адамс⁵ по поводу создания Вселенной, «это рассердило многих людей и повсеместно признавалось плохим ходом». И все потому, что программный код на Python 2.6 или 2.7 обычно не может исполняться интерпретатором Python 3, по крайней мере без небольшой модификации (хотя изменения, как правило, не слишком обременительны).

Всегда существует трение между неумолимой силой прогресса и обратной совместимостью. В этом случае команда разработчиков ядра языка Python решила, чтобы устранить некоторые несоответствия, в особенности в лежащем в основе программном интерфейсе C, требуется полный разрыв, и продвинула

¹ См. <http://jupyter.org/>.

² См. <http://pandas.pydata.org/>.

³ См. <http://scikit-learn.org/>.

⁴ См. <http://scikit-image.org/>.

⁵ См.: Адамс Д. Путеводитель автостопщика по Галактике (*Douglas Adams. The Hitchhiker's Guide to the Galaxy*. London: Pan Books, 1979).

язык в XXI век (Python 1.0 появился в 1994 г., более 20 лет назад, что в технологическом мире представляет собой целую жизнь).

Вот один из тех приемов, благодаря которым при переходе к версии 3 язык Python стал лучше:

```
print "Привет, Мир!" # инструкция print в Python 2
print("Привет, Мир!") # функция print в Python 3
```

Зачем создавать столько шума, чтобы добавить несколько круглых скобок! Все верно. Но если вместо этого вы хотите направить печать в другой *поток*, такой, например, как *стандартная ошибка*, т. е. в обычное место для отладочной информации?

```
print >>sys.stderr, "фатальная ошибка" # Python 2
print("фатальная ошибка", file=sys.stderr) # Python 3
```

Такое изменение, несомненно, выглядит целесообразнее. А что же происходит в Python версии 2? Авторы не знают, с полным на то правом.

Еще одно изменение состоит в том, как в Python 3 рассматривается деление целых чисел. Это делается точно так же, как большинство людей выполняет данную операцию. (Обратите внимание, что серия символов >>> говорит о том, что мы набираем программный код в интерактивной оболочке Python.)

```
# Python 2
>>> 5 / 2 2
# Python 3
>>> 5 / 2 2.5
```

Кроме того, нас также сильно порадовал новый оператор *умножения матриц* @, введенный в Python 3.5 в 2015 г. Обратитесь к главам 5 и 6, чтобы увидеть несколько практических примеров применения этого оператора!

Возможно, самым большим улучшением в Python 3 является поддержка Юникода, т. е. приема кодирования текста, который позволяет использовать не только английский алфавит, но и любой другой существующий в мире алфавит. Python 2 позволял вам определять строковое значение в кодировке Юникода следующим образом:

```
beta = u"β"
```

Но в Python 3 абсолютно все является Юникодом:

```
β = 0.5
print(2 * β)
```

```
1.0
```

Команда разработчиков ядра языка Python приняла абсолютно правильное решение, что в программном коде на Python имеет смысл в качестве объектов первого класса поддерживать символы всех естественных языков. Это в особенности актуально теперь, когда большинство новых программистов проживают не в англоязычных странах. Ради функциональной совместимости мы рекомендуем в большей части исходного кода по-прежнему использовать ан-

глийские символы. Эта возможность может пригодиться, например, в блокнотах Jupyter, утяжеленных математическими формулами.

☑ Наберите в терминале IPython или в блокноте Jupyter LaTeX-е имя символа, после чего нажмите клавишу Tab. В результате это имя будет расширено в символ Юникода. Например, `\beta<TAB>` станет β .

Обновление до Python 3 также разрушает большую часть существующего программного кода в версии 2.x и в некоторых случаях Python 3 выполняется медленнее, чем прежде. Несмотря на эти недостатки, мы рекомендуем всем пользователям как можно скорее обновиться до третьей версии (до 2020 г. Python 2.x теперь находится только в режиме обслуживания), поскольку большинство существовавших проблем будет решено вместе с совершенствованием версий 3.x. И действительно, в этой книге мы используем многие новые возможности Python 3.

В данной книге мы используем **Python 3.6**.

Дополнительные материалы для чтения, касающиеся перехода на версию 3, вы сможете найти на ресурсе Эда Шифилда Python-Future¹ и в справочнике книжного формата² Ника Коглана.

ЭКОСИСТЕМА И СООБЩЕСТВО SciPy

SciPy – это главная библиотека с довольно-таки большой функциональностью. Вместе с NumPy она является одним из уникальных приложений Python. Она положила начало огромному количеству связанных с ней библиотек, которые опираются на ее функционал. Со многими этими библиотеками вы будете работать на протяжении всей книги.

Создатели данных библиотек и многие их пользователи встречаются на многочисленных мероприятиях и конференциях по всему миру. Это такие мероприятия, как ежегодная конференция SciPy в Остин (США), EuroSciPy, SciPy Индия, PyData, и другие. Мы настоятельно вам рекомендуем посетить одну из них и встретиться с авторами лучшего научного программного обеспечения в мире Python. Если же у вас нет возможности попасть на эти конференции или же вы просто хотите прочувствовать их характер, воспользуйтесь Сетью³, где участники этих конференций публикуют свои дискуссии.

Бесплатное программное обеспечение и программное обеспечение с открытым исходным кодом (FOSS)

Сообщество SciPy приветствует разработку программного обеспечения с открытым исходным кодом. Исходный код почти всех библиотек SciPy находится

¹ См. <http://python-future.org/>.

² См. http://python-notes.curiousefficiency.org/en/latest/python3/questions_and_answers.html.

³ См. <https://www.youtube.com/user/EnthoughtMedia/playlists>.

в свободном доступе для чтения, правки и повторного использования компонентов любим, кто в этом заинтересован.

Если вы хотите, чтобы другие разработчики использовали ваш программный код, то один из лучших способов этого добиться состоит в том, чтобы сделать его свободным и открытым. Если вы используете программное обеспечение с закрытым исходным кодом и получаете не тот результат, который вы хотели достигнуть, то вам не повезло. Вы можете послать разработчику электронное сообщение и попросить, чтобы он добавил новый функционал (что часто не срабатывает!), либо написать новое программное обеспечение самостоятельно. Если исходный код находится в открытом доступе, то вы с легкостью можете добавить или изменить его функциональность, используя приемы, которые вы узнаете из этой книги.

Таким же образом, если вы находите в компоненте программного обеспечения системную ошибку, то наличие доступа к исходному коду может в значительной степени облегчить работу как пользователя, так и разработчика. Даже если вы не вполне разбираетесь в исходном коде, вы, как правило, продвинетесь в диагностике возникшей проблемы гораздо дальше и поможете разработчику с ее исправлением. Обычно это полезный познавательный опыт для всех!

Открытый исходный код, открытая наука

В научном программировании все вышеупомянутые сценарии чрезвычайно распространены и важны: научное программное обеспечение часто строится на предыдущей работе либо ее видоизменяет самым интересным образом. А вследствие высокого темпа научных публикаций и прогресса большое количество программного кода остается без тщательного тестирования перед его выпуском, что приводит к незначительным или системным ошибкам.

Еще одна веская причина, чтобы сделать исходный код открытым, состоит в том, чтобы способствовать развитию производимых исследований. Многие из нас по своему опыту знают, когда, читая действительно актуальную исследовательскую работу и затем скачивая исходный код, чтобы проверить его на своих собственных данных, мы обнаруживаем: исполняемый файл не скомпилирован для вашей операционной системы, невозможно разобраться, как его запустить, код имеет ошибки, отсутствуют важные компоненты. Или мы вообще получаем неожиданные результаты. Делая научное программное обеспечение открытым, мы не только улучшаем качество этого программного обеспечения, но и позволяем ясно увидеть, каким образом был написан код, какие допущения были приняты и жестко запрограммированы. Открытый исходный код помогает решать многие из этих проблем. Он также позволяет другим ученым опираться на исходный код своих коллег, способствуя новому сотрудничеству и ускоряя научный прогресс.

Лицензии на программное обеспечение с открытым исходным кодом

Если вы хотите, чтобы ваш программный код использовали другие, то вы *должны* его лицензировать. Если вы его не лицензируете, то по умолчанию он будет

закрытым. Даже если вы свой код опубликуете (например, разместив его в публичном хранилище GitHub), то без лицензии на программное обеспечение ваш программный код никто не имеет права использовать, править или распространять.

Выбирая среди многих вариантов лицензирования, сначала необходимо решить, что именно вы хотите позволить людям делать с вашим исходным кодом. Предоставить людям право продавать ваш исходный код для получения прибыли? Или право продавать программное обеспечение, в котором используется ваш исходный код? Или же вы хотите ограничить использование своего исходного кода только бесплатным программным обеспечением?

Есть две широкие категории FOSS-лицензий (лицензий на свободное и открытое программное обеспечение, Free and Open Source Software):

- разрешительная лицензия;
- свободная лицензия (Copy-left).

Разрешительная лицензия означает, что вы предоставляете любому пользователю право использовать, править и распространять ваш исходный код любым способом, который ему нравится, включая применение вашего исходного кода в качестве коммерческого программного обеспечения. В этой категории популярными вариантами являются лицензии MIT (программная лицензия Массачусетского технологического института) и BSD (программная лицензия университета Беркли). Сообщество SciPy приняло новую лицензию BSD (так называемую «Модифицированную BSD», или «3-пунктовую лицензию BSD»). Использование такой лицензии подразумевает получение помощи относительно исходного кода от огромного количества людей, включая тех, кто трудится в информационной индустрии и стартапах.

Свободные лицензии также позволяют другим разработчикам использовать, править и распространять ваш исходный код. Вместе с тем эти лицензии еще предписывают, что производный исходный код должен распространяться в соответствии со свободной лицензией. Таким образом свободные лицензии ограничивают то, что именно пользователи могут делать с этим исходным кодом.

Самой популярной свободной лицензией является Публичная лицензия GNU или GPL. Главный ее недостаток связан с использованием свободной лицензии и состоит в том, что часто ваш исходный код становится недоступным для любых потенциальных пользователей или участников из частного сектора. И среди них в будущем можете оказаться вы сами! Как результат этот факт может существенно уменьшить вашу пользовательскую базу и, следовательно, успех вашего программного обеспечения. В науке это может означать меньшее количество цитирования.

Для получения более подробной справки относительно выбора лицензии обратитесь на веб-сайт, посвященный выбору лицензии Choose a License¹. От-

¹ См. <http://choosealicense.com/>.

носителю лицензирования в научном контексте мы рекомендуем публикацию в блоге «The Whys and Hows of Licensing Scientific Code» (Вопросы «почему» и «как» относительно лицензирования научного программного кода) под авторством Джейка Вандерпласа (Jake VanderPlas), директора по исследованиям в области естествознания Вашингтонского университета и разносторонней суперзвезды SciPy. Собственно, здесь мы процитируем Джейка, чтобы убедительно довести до вас ключевые моменты лицензирования программного обеспечения:

...если вы извлечете из статьи всего три порции информации, то пусть они будут следующими:

1. Всегда лицензируйте свой код. Нелицензированный код является закрытым, поэтому любая открытая лицензия лучше, чем ничего (но см. п. 2).
2. Всегда используйте лицензию, совместимую с общей публичной лицензией (GPL). GPL-совместимые лицензии гарантируют вашему исходному коду широкую совместимость и включают GPL, новую BSD, MIT и другие (но см. п. 3).
3. Всегда используйте разрешительную лицензию в стиле BSD. Разрешительная лицензия, такая как новая лицензия BSD или лицензия MIT, более предпочтительна, по сравнению со свободной лицензией, такой как GPL или LGPL.

Весь исходный код в этой книге доступен в соответствии с 3-пунктовой лицензией BSD (3-clause BSD license). Там, где мы разместили фрагменты исходного кода других авторов. Этот исходный код обычно действует в соответствии с разрешительной открытой лицензией в той или иной форме (хотя не обязательно в соответствии с лицензией BSD).

Что касается вашего собственного исходного кода, то мы рекомендуем вам применять практику своего сообщества. В научном Python это означает 3-пунктовую лицензию BSD, в то время как в сообществе разработчиков на языке R, например, принята лицензия GPL.

GitHub: исходный код в социальном пространстве

Мы немного поговорили о распространении своего исходного кода в соответствии с лицензией на открытое программное обеспечение. И будем надеяться, что огромное число людей будут скачивать ваш исходный код, использовать его, исправлять ошибки и добавлять новые свойства. В связи с этим возникает вопрос: где разместить свой исходный код, чтобы люди смогли его найти? Как эти исправления ошибок и новые свойства вернуться в ваш исходный код? Каким образом отслеживать все вопросы и изменения? Можно представить, как все это совсем скоро выйдет из-под контроля.

Присоединяйтесь к GitHub.

GitHub¹ – это веб-сайт для размещения, совместного использования и разработки исходного кода. В его основе лежит программная система управления

¹ См. <https://github.com/>.

версиями Git¹. Обучению работе с GitHub посвящено несколько замечательных ресурсов, таких как «Введение в GitHub»² Питера Белла и Брента Бира. Подавляющее большинство проектов в экосистеме SciPy размещено на GitHub, поэтому, безусловно, стоит научиться его использовать!

Веб-сайт GitHub оказал значительное влияние на участие разработчиков и внесение ими открытого исходного кода. Это было достигнуто за счет предоставления пользователям возможности публиковать исходный код и свободно сотрудничать. Любой участник может зайти и создать копию (так называемое ответвление *fork*) исходного кода и отредактировать, как его душе будет угодно. Он в конечном счете может внести эти изменения назад в оригинальный исходный код, создав *запрос на включение внесенных изменений*. Существует целый ряд хороших возможностей, таких как управление текущими вопросами и запросами на изменение, а также возможность определять, кто непосредственно может редактировать ваш исходный код. Вы даже можете отслеживать правки участников и получать другую детальную статистику. Помимо вышеперечисленного, существует целый ряд других замечательных особенностей GitHub. Однако мы предоставим вам самим возможность обнаружить многие из них самостоятельно. С некоторыми из них вы познакомитесь, когда будете читать последующие главы. В сущности, GitHub демократизировал разработку программного обеспечения (рис. П.1) и существенно снизил барьер доступа к нему.

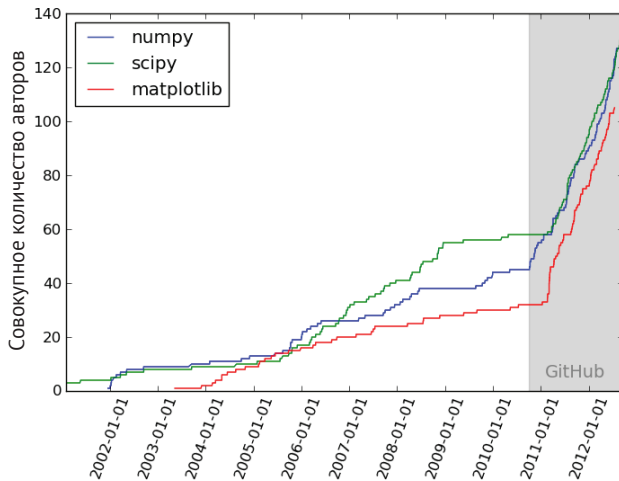


Рис. П.1 ❖ Влияние GitHub
(используется с разрешения автора, Джейка Вандерпласа)

¹ См. <http://git-scm.com/>.

² См. <http://shop.oreilly.com/product/0636920033059.do>.

Оставьте свой след в экосистеме SciPy

По мере накопления опыта работы с SciPy и после того, как вы начнете использовать эту библиотеку в вашей исследовательской работе, может оказаться, что тому или иному пакету не хватает функционала, в котором вы нуждаетесь. Либо вы посчитаете, что можете сделать что-то эффективнее. А возможно, найдете ошибку. Когда вы дойдете до этого, значит, пришла пора начать принимать участие в экосистеме SciPy.

Мы настоятельно рекомендуем вам попробовать. Сообщество существует, потому что люди готовы делиться своим исходным кодом и совершенствовать существующий исходный код. И если каждый из нас внесет свой небольшой вклад, то вместе мы добьемся многого. Но помимо любых альтруистических причин содействия сообществу, существует несколько вполне практических личных выгод. Сотрудничая с сообществом, вы станете более профессиональным программистом. Любой вносимый вами исходный код будет рассмотрен другими участниками, и вы получите отзыв в виде замечаний и комментариев. В качестве побочного эффекта вы научитесь использовать Git и GitHub, которые сами по себе являются очень полезными инструментами для технической поддержки и совместного использования вашего собственного исходного кода. Вы даже можете обнаружить, что взаимодействие с сообществом SciPy предоставляет вам более широкую научную сеть и удивительные возможности карьерного роста.

Мы хотим, чтобы вы задумались над тем, чтобы стать больше, чем просто пользователем SciPy. Присоединившись к сообществу, ваша работа сделает его лучше для всех научных программистов.

Капля эксцентричности Python

Если вы обеспокоены тем, что сообщество SciPy может стать для вновь прибывшего участника вызывающим благоговение местом, то следует напомнить, что это сообщество состоит из таких же людей, как и вы, ученых, которые, как правило, имеют прекрасное чувство юмора.

В стране Python вы неизбежно найдете ссылки на шоу «Монти Пайтон»¹. Пакет *Airspeed Velocity*² измеряет быстродействие вашего программного обеспечения (подробнее о нем позднее) и ссылается на строчку из «Монти Пайтона и Священного Грааля» «какова воздушная скорость полета необремененной ласточки?».

А вот еще один забавно названный пакет – «Six». Он позволяет использовать пакеты Python 2 из Python 3. В имени пакета обыгрывается слово «six» (шесть),

¹ Монти Пайтон (англ. Monty Python) – комик-группа из Великобритании, состоявшая из шести человек. Благодаря своему новаторскому, абсурдистскому юмору участники «Монти Пайтон» находятся в числе самых влиятельных комиков всех времен. Группа известна во многом благодаря юмористическому телешоу «Летающий цирк Монти Пайтона». – *Прим. перев.*

² См. <http://spacetelescope.github.io/asv/using.html>.

а сам пакет позволяет использовать синтаксис Python 3 в Python 2 с новозеландским акцентом. Синтаксис `Sux` уменьшает разочарование от использования пакетов, предназначенных только для Python 2, после того как вы перешли на Python 3:

```
import sux
p = sux.to_use('my_py2_package')
```

В целом имена библиотек Python могут вызывать бурное веселье, и мы надеемся, что вы хорошо проведете время, когда будете придумывать свое собственное!

Получение помощи

Когда мы попадаем в тупик, наш первый шаг – поискать в Интернете подсказку, которая поможет найти решение текущей задачи, либо информацию об ошибке, сообщение о которой мы получили. Эти поиски обычно приводят нас на Stack Overflow¹, превосходный вопросно-ответный сайт для программистов. Если вы с первого раза не нашли того, что ищете, то попробуйте обобщить критерии поиска, чтобы найти кого-то, кто имеет схожие проблемы.

Иногда вы можете оказаться фактически первым человеком, который задаст этот конкретный вопрос (это наиболее вероятно в тех случаях, когда вы используете совершенно новый пакет). Но не отчаивайтесь! Не все потеряно! Как было отмечено выше, сообщество SciPy дружелюбно и разбросано в различных частях межсетей. Ваш следующий шаг состоит в том, чтобы обратиться в поисковик с запросом «<имя библиотеки> список рассылки» и найти список рассылки с адресами, куда можно обратиться за помощью. Авторы библиотек и компетентные пользователи регулярно их читают и очень радушны ко вновь прибывшим участникам. Обратите внимание, что общепринятое правило поведения в сообществе требует, чтобы, перед тем как отправлять свои вопросы, вы *подписались* на список. Если вы не подпишетесь, это будет означать, что перед регистрацией вашего электронного адреса в списке кто-то должен будет вручную проверить, что данный адрес не является спамным. Присоединение к очередному списку рассылки может показаться ненужным, но мы настоятельно рекомендуем поступать именно так: это как раз то место, где надо учиться!

Инсталляция языка Python

В настоящей книге мы исходим из того, что у вас уже установлен Python 3.6 (либо его более поздняя версия) и все необходимые пакеты SciPy. Мы перечислим все использованные нами необходимые библиотеки и их версии в файле `environment.yml`, прилагаемом вместе с файлами данных и исходными кодами

¹ См. <http://stackoverflow.com/>.

к этой книге. Самый легкий способ получить все необходимые компоненты – установить `conda`¹, инструмент для управления средой разработки на Python. Затем вы можете передать файл `environment.yml` в `conda`, чтобы за один шаг установить правильные версии всех необходимых пакетов.

```
conda env create --name elegant-scipy -f путь/к/environment.yml
source activate elegant-scipy
```

За дополнительной информацией обратитесь к хранилищу книги на GitHub².

Доступ к книжным материалам

Весь исходный код и данные, приводимые в этой книге, доступны в нашем хранилище на GitHub. В файле README хранилища вы найдете инструкции по созданию блокнотов Jupyter из специальных исходных файлов с облегченной разметкой markdown. После этого блокноты Jupyter можно запускать в интерактивном режиме, используя данные, которые включены в хранилище.

НАЧИНАЕМ

Мы подобрали самый элегантный исходный код, который предлагает сообщество SciPy. Мы также займемся исследованием нескольких реальных научных задач, решаемых при помощи SciPy. Эта книга еще дает возможность мимолетно взглянуть на радушное и готовое к сотрудничеству научное программистское сообщество, которое надеется, что вы присоединитесь.

Добро пожаловать в элегантный SciPy.

УСЛОВНЫЕ ОБОЗНАЧЕНИЯ, ПРИНЯТЫЕ В КНИГЕ

В книге используются следующие типографские условные обозначения.

Курсивный шрифт

Указывает новые термины, URL-адреса, адреса электронной почты, имена и расширения файлов.

Моноширинный шрифт

Используется для листингов программ, а также внутри абзацев для отсылки на элементы программ, такие как переменные или имена функций, базы данных, типы данных, переменные окружающей среды, операторы и ключевые слова.

Жирный моноширинный шрифт




Показывает команды либо другой текст, который должен быть напечатан самим пользователем.

¹ См. <http://conda.pydata.org/miniconda.html>.

² См. <https://github.com/elegant-scipy/elegant-scipy>.

Курсивный моноширинный шрифт

Показывает текст, который должен быть заменен на предоставленные пользователем значения либо на значения, определяемые контекстом.

-  Данный элемент обозначает общее замечание.
-  Данный элемент обозначает подсказку или совет.
-  Данный элемент обозначает предупреждение или предостережение.

ИСПОЛЬЗОВАНИЕ ЦВЕТА

В некоторых приводимых в книге примерах используются различные цвета, которые не видны в печатной версии этой книги. Читателям печатной версии книги рекомендуется обратиться к исходным блокнотам на <http://elegant-scipy.org/>.

ИСПОЛЬЗОВАНИЕ ПРИМЕРОВ ПРОГРАММ

Дополнительный материал (примеры программного кода, упражнения и т. д.) доступен для скачивания с <https://github.com/elegant-scipy/elegant-scipy>.

Эта книга предназначена, чтобы помочь вам в решении своих задач. В целом, если код примеров предлагается вместе с книгой, вы можете использовать его в своих программах и документации. Вам не нужно связываться с нами с просьбой о разрешении, если вы не воспроизводите значительную часть кода. Например, написание программы, которая использует несколько фрагментов кода из данной книги, официального разрешения не требует. Продажа либо распространение компакт-диска с примерами из книг издательства O'Reilly требует официального разрешения. Ответ на вопрос цитированием данной книги и приведение выдержек кода примеров в качестве цитат разрешения не требует. Включение значительного количества кода примеров из данной книги в документацию на ваш продукт требует официального разрешения.

Мы ценим, но не требуем атрибуции. Атрибуция обычно включает в себя титул, автора, издателя и ISBN. Например: «Нуньес-Иглесиас Х., Уолт ван дер Ш., Дэшноту Х. Элегантный SciPy. O'Reilly. ISBN 978-1-491-92287-3».

Если вы считаете, что применяете примеры кода, выходя за рамки их справедливого использования или разрешения, выданного выше, то обращайтесь к нам по адресу permissions@oreilly.com.

БЛАГОДАРНОСТИ

Выражаем признательность многим и многим людям, которые внесли существенный вклад в эту книгу. Она бы не появилась без вашей помощи.

Конец ознакомительного фрагмента.

Приобрести книгу можно

в интернет-магазине

«Электронный универс»

e-Univers.ru