
Grundlagen von NumPy: Arrays und vektorisierte Berechnung

NumPy, kurz für Numerical Python, ist eines der wichtigsten Pakete für numerische Berechnungen in Python. Die meisten wissenschaftlichen Rechenpakete nutzen die Array-Objekte von NumPy als *Lingua Franca* für den Datenaustausch.

Hier sind einige der Dinge, die Sie in NumPy finden werden:

- ndarray, ein effizientes mehrdimensionales Array, das schnelle Array-orientierte Arithmetikoperationen und flexible *Broadcasting*-Fähigkeiten mitbringt.
- Mathematische Funktionen für schnelle Operationen auf ganzen Daten-Arrays, ohne dass man Schleifen schreiben muss.
- Tools zum Lesen und Schreiben von Array-Daten auf die Festplatte und zum Arbeiten mit Memory-mapped Dateien.
- Lineare Algebra, Generierung von Zufallszahlen und Fourier-Transformation.
- Eine C-API zum Anbinden von NumPy an Bibliotheken, die in C, C++ oder Fortran geschrieben sind.

Wegen der leicht zu benutzenden C-API von NumPy ist es einfach, Daten an externe Bibliotheken weiterzugeben, die in einer niedrigen Programmiersprache geschrieben sind. Ebenso einfach können externe Bibliotheken Daten als NumPy-Arrays an Python zurückliefern. Diese Eigenschaft macht Python zur Sprache der Wahl, wenn man Wrapper für vorhandenen C-/C++-/Fortran-Code schreiben und diesem eine dynamische und leicht zu verwendende Schnittstelle geben möchte.

NumPy selbst bietet zwar keine wissenschaftliche oder Modellierungsfunktionalität, dennoch hilft Ihnen ein Grundverständnis von NumPy-Arrays und Array-orientierten Berechnungen bei der effektiveren Benutzung von Tools mit Array-orientierter Semantik wie etwa pandas. Da NumPy ein so umfassendes Thema ist, werde ich viele komplexere NumPy-Funktionen wie etwa Broadcasting später ausführlicher behandeln (siehe Anhang A).

Für die meisten Datenanalyseanwendungen konzentriere ich mich vor allem auf folgende Bereiche:

- Schnelle vektorisierte Array-Operationen zum Bereinigen von Daten, Bilden und Filtern von Teilmengen sowie für Transformationen und andere Arten von Berechnungen.
- Gebräuchliche Array-Algorithmen wie Sortieren, Eindeutigkeit und Mengenoperationen.
- Effiziente beschreibende Statistik und das Aggregieren/Zusammenfassen von Daten.
- Datenausrichtung und relationale Datenmanipulationen zum Mischen und Verbinden heterogener Datenmengen.
- Ausdrücken logischer Bedingungen als Array-Ausdruck statt als Schleifen mit `if-elif-else`-Verzweigungen.
- Gruppenweise Datenmanipulation (Aggregation, Transformation, Funktionsanwendung).

Auch wenn NumPy die rechnerische Grundlage für die allgemeine numerische Datenverarbeitung bietet, werden viele Leser pandas als Basis für die meisten Arten von Statistiken oder Analysen, speziell an Tabellendaten, verwenden wollen. pandas stellt darüber hinaus eine etwas speziellere Funktionalität bereit, wie die Manipulation von Zeitreihen, die es in NumPy nicht gibt.



Die Wurzeln des Array-orientierten Rechnens in Python lassen sich bis 1995 zurückverfolgen, als Jim Hugunin die Numeric-Bibliothek geschaffen hat. Im Laufe der nächsten zehn Jahre begannen viele wissenschaftliche Gemeinschaften mit der Array-Programmierung in Python, allerdings kam es mit Beginn der 2000er-Jahre zu einer Zersplitterung des Bibliotheksökosystems. 2005 schaffte es Travis Oliphant, aus den damaligen Numeric- und Numarray-Projekten das NumPy-Projekt zusammenzuschmieden, um der Community ein einziges Array-Rechenframework zur Verfügung zu stellen.

Einer der Gründe dafür, dass NumPy für numerische Berechnungen in Python so wichtig ist, besteht in seiner Effizienz bei großen Daten-Arrays. Das hat verschiedene Ursachen:

- NumPy speichert Daten intern in einem zusammenhängenden Speicherblock unabhängig von anderen eingebauten Python-Objekten. Die NumPy-Algorithmenbibliothek, die in C geschrieben ist, kann ohne irgendwelche Typprüfungen oder anderen Aufwand auf diesem Speicher arbeiten. NumPy-Arrays verwenden außerdem viel weniger Speicher als eingebaute Python-Sequenzen.
- NumPy-Operationen führen komplexe Berechnungen auf ganzen Arrays aus, ohne dass sie dazu Python-`for`-Schleifen brauchen.

Damit Sie eine Vorstellung von dem Leistungsunterschied bekommen, betrachten Sie ein NumPy-Array mit einer Million Integer-Werten und die entsprechende Python-Liste:

```
In [7]: import numpy as np
```

```
In [8]: my_arr = np.arange(1000000)
```

```
In [9]: my_list = list(range(1000000))
```

Multiplizieren wir nun jede Sequenz mit 2:

```
In [10]: %time for _ in range(10): my_arr2 = my_arr * 2  
CPU times: user 20 ms, sys: 50 ms, total: 70 ms  
Wall time: 72.4 ms
```

```
In [11]: %time for _ in range(10): my_list2 = [x * 2 for x in my_list]  
CPU times: user 760 ms, sys: 290 ms, total: 1.05 s  
Wall time: 1.05 s
```

NumPy-basierte Algorithmen sind im Allgemeinen 10- bis 100-mal schneller (oder mehr) als ihre reinen Python-Gegenstücke und nutzen dabei deutlich weniger Speicher.

4.1 Das ndarray von NumPy: ein mehrdimensionales Array-Objekt

Eines der zentralen Merkmale von NumPy ist sein N-dimensionales Array-Objekt oder ndarray, ein schneller, flexibler Container für große Datenmengen in Python. Arrays erlauben Ihnen, mathematische Operationen auf ganzen Datenblöcken durchzuführen, wobei die Syntax den äquivalenten Operationen zwischen skalaren Elementen ähnlich ist.

Damit Sie einen Vorgeschmack darauf bekommen, wie NumPy Stapelberechnungen mit einer ähnlichen Syntax durchführt wie bei den Skalaren auf eingebauten Python-Objekten, importiere ich zuerst NumPy und generiere ein kleines Array von Zufallsdaten:

```
In [12]: import numpy as np
```

```
# Generate some random data
```

```
In [13]: data = np.random.randn(2, 3)
```

```
In [14]: data
```

```
Out[14]:  
array([[ -0.2047,  0.4789, -0.5194],  
       [-0.5557,  1.9658,  1.3934]])
```

Anschließend schreibe ich mathematische Operationen mit `data`:

```
In [15]: data * 10
Out[15]:
array([[ -2.0471,   4.7894,  -5.1944],
       [-5.5573,  19.6578,  13.9341]])
```

```
In [16]: data + data
Out[16]:
array([[ -0.4094,   0.9579,  -1.0389],
       [-1.1115,   3.9316,   2.7868]])
```

Im ersten Beispiel wurden alle Elemente mit 10 multipliziert, im zweiten wurden die korrespondierenden Werte in jeder »Zelle« des Arrays miteinander addiert.



In diesem Kapitel und auch im Rest des Buchs verfare ich immer nach der Konvention `import numpy as np`. Sie dürfen natürlich gern `from numpy import *` in Ihrem Code schreiben, um das ewige `np.` zu vermeiden, allerdings rate ich dringend davon ab, sich das anzugewöhnen. Der `numpy`-Namensraum ist groß und enthält eine Reihe von Funktionen, deren Namen mit den eingebauten Python-Funktionen in Konflikt geraten (wie `min` und `max`).

Ein `ndarray` ist ein generischer, mehrdimensionaler Container für homogene Daten, d. h., alle Elemente müssen vom selben Typ sein. Jedes Array hat einen `shape`, ein Tupel, das die Größe jeder Dimension angibt, und einen `dtype`, ein Objekt, das den Datentyp des Arrays beschreibt:

```
In [17]: data.shape
Out[17]: (2, 3)
```

```
In [18]: data.dtype
Out[18]: dtype('float64')
```

Dieses Kapitel führt Sie in die Grundlagen der Benutzung von NumPy-Arrays ein und sollte Sie ausreichend auf den Rest des Buchs vorbereiten. Es ist für viele datenanalytische Anwendungen zwar nicht nötig, ein tiefes Verständnis von NumPy zu haben, allerdings hilft es Ihnen auf dem Weg zum Guru des wissenschaftlichen Python, wenn Sie sich mit der Array-orientierten Programmierung und Denkweise vertraut machen.



Die Begriffe »Array«, »NumPy-Array« und »ndarray« beziehen sich mit wenigen Ausnahmen immer auf dasselbe: das `ndarray`-Objekt.

ndarrays erzeugen

Am einfachsten erzeugen Sie ein Array mit der `array`-Funktion. Diese akzeptiert jedes sequenzartige Objekt (einschließlich anderer Arrays) und erzeugt ein neues

NumPy-Array, das die übergebenen Daten enthält. So ist zum Beispiel eine Liste ein guter Kandidat für eine Konvertierung:

```
In [19]: data1 = [6, 7.5, 8, 0, 1]
In [20]: arr1 = np.array(data1)
In [21]: arr1
Out[21]: array([ 6. ,  7.5,  8. ,  0. ,  1. ])
```

Verschachtelte Sequenzen, wie eine Liste aus Listen gleicher Länge, werden in ein mehrdimensionales Array umgewandelt:

```
In [22]: data2 = [[1, 2, 3, 4], [5, 6, 7, 8]]
In [23]: arr2 = np.array(data2)
In [24]: arr2
Out[24]:
array([[1, 2, 3, 4],
       [5, 6, 7, 8]])
```

Da `data2` eine Liste von Listen war, hat das NumPy-Array `arr2` zwei Dimensionen mit einer Form, die aus den Daten abgeleitet wurde. Wir können das bestätigen, wenn wir die Attribute `ndim` und `shape` untersuchen:

```
In [25]: arr2.ndim
Out[25]: 2
In [26]: arr2.shape
Out[26]: (2, 4)
```

Wenn nichts explizit angegeben ist (mehr dazu später), versucht `np.array`, einen guten Datentyp für das Array abzuleiten, das es erzeugt. Der Datentyp wird in einem speziellen `dtype`-Metadatenobjekt gespeichert. So haben wir etwa in den vorherigen zwei Beispielen Folgendes:

```
In [27]: arr1.dtype
Out[27]: dtype('float64')
In [28]: arr2.dtype
Out[28]: dtype('int64')
```

Neben `np.array` gibt es eine Reihe weiterer Funktionen zum Anlegen neuer Arrays. `zeros` und `ones` erzeugen zum Beispiel Arrays aus Nullen bzw. Einsen bei vorgegebener Länge oder Form. `empty` erzeugt ein Array, ohne dessen Werte auf einen bestimmten Wert zu initialisieren. Um mit diesen Methoden ein höherdimensionales Array zu erzeugen, übergeben Sie für die Form ein Tupel:

```
In [29]: np.zeros(10)
Out[29]: array([ 0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.])
In [30]: np.zeros((3, 6))
Out[30]:
array([[ 0.,  0.,  0.,  0.,  0.,  0.],
       [ 0.,  0.,  0.,  0.,  0.,  0.]])
```

```
[ 0., 0., 0., 0., 0., 0.]])
```

```
In [31]: np.empty((2, 3, 2))  
Out[31]:  
array([[ [ 0., 0.],  
        [ 0., 0.],  
        [ 0., 0.]],  
       [[ [ 0., 0.],  
        [ 0., 0.],  
        [ 0., 0.]])
```



Man darf nicht davon ausgehen, dass `np.empty` ein Array zurückliefert, das nur Nullen enthält. Manchmal steht auch einfach Müll drin.

`arange` ist eine Version der eingebauten Python-Funktion `range` für Arrays:

```
In [32]: np.arange(15)  
Out[32]: array([ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14])
```

In Tabelle 4-1 finden Sie eine kurze Liste der Standardfunktionen zum Erzeugen von Arrays. Da sich NumPy auf das numerische Rechnen konzentriert, ist der Datentyp in vielen Fällen `float64` (Gleitkommazahl), sofern nichts angegeben ist.

Tabelle 4-1: Funktionen zum Erzeugen von Arrays

Funktion	Beschreibung
<code>array</code>	Wandelt Eingabedaten (Liste, Tupel, Array oder ein anderer Sequenztyp) in ein <code>ndarray</code> um, indem entweder ein <code>dtype</code> abgeleitet oder explizit ein <code>dtype</code> angegeben wird; kopiert standardmäßig die Eingabedaten.
<code>asarray</code>	Wandelt Eingabedaten in ein <code>ndarray</code> um, kopiert sie aber nicht, falls die Eingabe bereits ein <code>ndarray</code> ist.
<code>arange</code>	Wie die eingebaute <code>range</code> -Funktion, liefert aber ein <code>ndarray</code> statt einer Liste zurück.
<code>ones</code> , <code>ones_like</code>	Erzeugt ein Array nur aus Einsen mit der angegebenen Form und <code>dtype</code> ; <code>ones_like</code> nimmt ein Array entgegen und erzeugt ein Einsen-Array derselben Form und <code>dtype</code> .
<code>zeros</code> , <code>zeros_like</code>	Wie <code>ones</code> und <code>ones_like</code> , erzeugt aber stattdessen Arrays aus Nullen.
<code>empty</code> , <code>empty_like</code>	Erzeugt neue Arrays, indem es neuen Speicher belegt, diesen aber anders als <code>ones</code> und <code>zeros</code> nicht mit Werten füllt.
<code>full</code> , <code>full_like</code>	Erzeugt ein Array der angegebenen Form und <code>dtype</code> mit Werten, die auf den angegebenen »Füllwert« gesetzt sind, <code>full_like</code> nimmt ein anderes Array entgegen und erzeugt ein gefülltes Array derselben Form und <code>dtype</code> .
<code>eye</code> , <code>identity</code>	Erzeugt eine quadratische $N \times N$ -Einheitsmatrix (Einsen auf der Diagonalen und ansonsten Nullen).

Datentypen für `ndarrays`

Der *Datentyp* oder `dtype` ist ein besonderes Objekt, das die Informationen (oder *Metadaten*, also Daten über Daten) enthält, die das `ndarray` benötigt, um einen Auszug aus dem Speicher als eine bestimmte Art von Daten zu identifizieren:

```
In [33]: arr1 = np.array([1, 2, 3], dtype=np.float64)
```

```
In [34]: arr2 = np.array([1, 2, 3], dtype=np.int32)
```

```
In [35]: arr1.dtype  
Out[35]: dtype('float64')
```

```
In [36]: arr2.dtype  
Out[36]: dtype('int32')
```

dtypes sind einer der Gründe für NumPys Flexibilität bei der Interaktion mit Daten aus anderen Systemen. In den meisten Fällen erlauben sie eine direkte Zuordnung zur zugrunde liegenden Festplatten- oder Speicherrepräsentation, was es leicht macht, binäre Datenströme zu lesen und Anbindungen an Code in maschinennahen Sprachen wie C oder Fortran vorzunehmen. Die numerischen dtypes werden auf die gleiche Weise benannt: ein Typname, wie float oder int, gefolgt von einer Zahl, die die Anzahl der Bits pro Element angibt. Ein normaler Gleitkommawert mit doppelter Genauigkeit (im Prinzip das Python-Objekt float) nimmt 8 Bytes oder 64 Bits ein. Deshalb wird dieser Typ in NumPy als float64 bezeichnet. Tabelle 4-2 zeigt eine vollständige Liste der von NumPy unterstützten Datentypen.



Sie müssen sich die NumPy-dtypes nicht merken, speziell als Anfänger nicht. Oft reicht es, sich um die allgemeine Art der Daten zu kümmern, mit denen Sie es zu tun haben, also Gleitkomma, Komplex, Integer, Boolean, String oder allgemeines Python-Objekt. Wenn Sie mehr Kontrolle darüber haben müssen, wie die Daten im Speicher und auf der Festplatte abgelegt werden, vor allem bei großen Datenmengen, ist es gut zu wissen, dass Sie den Speichertyp kontrollieren können.

Tabelle 4-2: NumPy-Datentypen

Typ	Typcode	Beschreibung
int8, uint8	i1, u1	Vorzeichenbehaftete und vorzeichenlose 8-Bit-(1-Byte-)Integer-Typen.
int16, uint16	i2, u2	Vorzeichenbehaftete und vorzeichenlose 16-Bit-Integer-Typen.
int32, uint32	i4, u4	Vorzeichenbehaftete und vorzeichenlose 32-Bit-Integer-Typen.
int64, uint64	i8, u8	Vorzeichenbehaftete und vorzeichenlose 64-Bit-Integer-Typen.
float16	f2	Gleitkommazahlen mit halber Genauigkeit.
float32	f4 oder f	Gleitkommazahlen mit einfacher Genauigkeit; kompatibel mit einem C float.
float64	f8 oder d	Gleitkommazahlen mit doppelter Genauigkeit; kompatibel mit einem C double und dem Python-Objekt float.
float128	f16 oder g	Gleitkommazahlen mit erweiterter Genauigkeit.
complex64, complex128, complex256	c8, c16, c32	Komplexe Zahlen, repräsentiert durch zwei Gleitkommazahlen mit 32, 64 bzw. 128 Bits
bool	?	Boolescher Typ, speichert die Werte True und False.
object	0	Python-Objektyp; Wert kann ein beliebiges Python-Objekt sein.

Tabelle 4-2: NumPy-Datentypen (Fortsetzung)

Typ	Typcode	Beschreibung
string_	S	ASCII-String-Typ mit fester Länge (1 Byte pro Zeichen); um zum Beispiel einen String dtype der Länge 10 zu erzeugen, nutzen Sie 'S10'.
unicode_	U	Unicode-Typ fester Länge (die Anzahl der Bytes ist plattformabhängig); die Semantik der Spezifikation ist identisch mit string_ (z.B. 'U10').

Sie können ein Array ausdrücklich von einem dtype in einen anderen umwandeln. Für diesen auch *Casting* genannten Vorgang nutzen Sie die ndarray-Methode `astype`:

```
In [37]: arr = np.array([1, 2, 3, 4, 5])

In [38]: arr.dtype
Out[38]: dtype('int64')

In [39]: float_arr = arr.astype(np.float64)

In [40]: float_arr.dtype
Out[40]: dtype('float64')
```

In diesem Beispiel wurden Integer-Werte in Gleitkommazahlen umgewandelt. Beim Umwandeln von Gleitkommazahlen in Integer wird der Nachkommateil abgeschnitten:

```
In [41]: arr = np.array([3.7, -1.2, -2.6, 0.5, 12.9, 10.1])

In [42]: arr
Out[42]: array([ 3.7, -1.2, -2.6,  0.5, 12.9, 10.1])

In [43]: arr.astype(np.int32)
Out[43]: array([ 3, -1, -2,  0, 12, 10], dtype=int32)
```

Falls Sie ein Array aus Strings haben, die Zahlen darstellen, können Sie sie mit `astype` in eine numerische Form konvertieren:

```
In [44]: numeric_strings = np.array(['1.25', '-9.6', '42'], dtype=np.string_)

In [45]: numeric_strings.astype(float)
Out[45]: array([ 1.25, -9.6 , 42.  ])
```



Seien Sie vorsichtig bei der Benutzung des Typs `numpy.string_`. String-Daten in NumPy haben nämlich eine feste Größe, und Eingaben könnten daher ohne Warnung abgeschnitten werden. pandas verhält sich bei nicht numerischen Daten deutlich intuitiver.

Sollte das Casting aus irgendeinem Grund fehlschlagen (weil etwa ein String nicht in float64 konvertiert werden kann), wird ein `ValueError` ausgelöst. Ich war hier ein bisschen faul und habe `float` statt `np.float64` geschrieben; NumPy setzt clevererweise für die Python-Typen seine eigenen äquivalenten dtypes ein.

Sie können auch das dtype-Attribut eines anderen Arrays benutzen:

```
In [46]: int_array = np.arange(10)
```

```
In [47]: calibers = np.array([.22, .270, .357, .380, .44, .50], dtype=np.float64)
```

```
In [48]: int_array.astype(calibers.dtype)
```

```
Out[48]: array([ 0.,  1.,  2.,  3.,  4.,  5.,  6.,  7.,  8.,  9.])
```

Es gibt Kurzformen für die Codes, mit denen Sie sich ebenfalls auf einen dtype beziehen können:

```
In [49]: empty_uint32 = np.empty(8, dtype='u4')
```

```
In [50]: empty_uint32
```

```
Out[50]:
```

```
array([          0, 1075314688,          0, 1075707904,          0,
        1075838976,          0, 1072693248], dtype=uint32)
```



Der Aufruf von `astype` erzeugt *immer* ein neues Array (eine Kopie der Daten), selbst wenn der neue dtype identisch mit dem alten dtype ist.

Rechnen mit NumPy-Arrays

Arrays sind wichtig, weil sie Ihnen erlauben, viele Operationen auf Daten auszuführen, ohne dass Sie `for`-Schleifen schreiben müssen. NumPy-Anwender bezeichnen dies als *Vektorisierung*. Jede arithmetische Operation zwischen Arrays gleicher Größe führt ihre Arbeit elementweise durch:

```
In [51]: arr = np.array([[1., 2., 3.], [4., 5., 6.]])
```

```
In [52]: arr
```

```
Out[52]:
```

```
array([[ 1.,  2.,  3.],
        [ 4.,  5.,  6.]])
```

```
In [53]: arr * arr
```

```
Out[53]:
```

```
array([[ 1.,  4.,  9.],
        [16., 25., 36.]])
```

```
In [54]: arr - arr
```

```
Out[54]:
```

```
array([[ 0.,  0.,  0.],
        [ 0.,  0.,  0.]])
```

Arithmetische Operationen mit Skalaren propagieren das skalare Argument zu jedem Element in dem Array:

```
In [55]: 1 / arr
```

```
Out[55]:
```

```
array([[ 1.    ,  0.5   ,  0.3333],
```

```
[ 0.25 , 0.2 , 0.1667]])
```

```
In [56]: arr ** 0.5  
Out[56]:  
array([[ 1.    ,  1.4142,  1.7321],  
       [ 2.    ,  2.2361,  2.4495]])
```

Vergleiche zwischen Arrays derselben Größe ergeben boolesche Arrays:

```
In [57]: arr2 = np.array([[0., 4., 1.], [7., 2., 12.]])
```

```
In [58]: arr2  
Out[58]:  
array([[ 0.,  4.,  1.],  
       [ 7.,  2., 12.]])
```

```
In [59]: arr2 > arr  
Out[59]:  
array([[False,  True, False],  
       [ True, False,  True]], dtype=bool)
```

Operationen zwischen Arrays unterschiedlicher Größe werden als *Broadcasting* bezeichnet. Wir gehen in Anhang A näher darauf ein. Für den Großteil dieses Buchs ist es nicht notwendig, ein tiefergehendes Verständnis zum Broadcasting zu entwickeln.

Einfaches Indizieren und Slicing

Das Indizieren von Arrays in NumPy ist ein weites Feld, weil es sehr viele Möglichkeiten gibt, eine Teilmenge Ihrer Daten oder einzelne Elemente auszuwählen. Eindimensionale Arrays sind einfach; oberflächlich betrachtet, verhalten sie sich ähnlich wie Python-Listen:

```
In [60]: arr = np.arange(10)  
  
In [61]: arr  
Out[61]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])  
  
In [62]: arr[5]  
Out[62]: 5  
  
In [63]: arr[5:8]  
Out[63]: array([5, 6, 7])  
  
In [64]: arr[5:8] = 12  
  
In [65]: arr  
Out[65]: array([ 0,  1,  2,  3,  4, 12, 12, 12,  8,  9])
```

Wie Sie sehen können, wird ein skalarer Wert, wenn Sie ihn einem Teilbereich wie `arr[5:8] = 12` zuweisen, an die gesamte Auswahl propagiert (*Broadcasting*). Ein erster wichtiger Unterschied zu Python's eingebauten Listen besteht darin, dass Teilbereiche (Slices) von Arrays sogenannte *Views* auf das ursprüngliche Array

sind. Das bedeutet, dass die Daten nicht kopiert werden und sich alle Modifikationen an dem View im Quell-Array niederschlagen.

Als Beispiel erzeuge ich zuerst ein Slice von arr:

```
In [66]: arr_slice = arr[5:8]
```

```
In [67]: arr_slice  
Out[67]: array([12, 12, 12])
```

Wenn ich nun die Werte in arr_slice ändere, spiegeln sich die Änderungen im Original-Array arr wider:

```
In [68]: arr_slice[1] = 12345
```

```
In [69]: arr  
Out[69]: array([ 0,  1,  2,  3,  4, 12, 12345, 12,  8,  
                9])
```

Das »leere« Slice [:] wird allen Werten in einem Array zugewiesen:

```
In [70]: arr_slice[:] = 64
```

```
In [71]: arr  
Out[71]: array([ 0,  1,  2,  3,  4, 64, 64, 64,  8,  9])
```

Sollten Sie neu bei NumPy sein, mag Sie das überraschen, vor allem wenn Sie andere Array-Programmiersprachen benutzt haben, die viel bereitwilliger Daten kopieren. Da NumPy mit sehr großen Arrays funktionieren soll, können Sie sich vorstellen, dass es zu Leistungs- und Speicherproblemen kommen könnte, wenn NumPy darauf bestehen würde, die Daten immer zu kopieren.



Falls Sie statt eines Views tatsächlich die Kopie eines Slice eines ndarray haben wollen, müssen Sie das Array explizit kopieren – zum Beispiel so: `arr[5:8].copy()`.

Bei höherdimensionalen Arrays haben Sie deutlich mehr Optionen. In einem zweidimensionalen Array sind die Elemente an jedem Index keine Skalare mehr, sondern eindimensionale Arrays:

```
In [72]: arr2d = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
```

```
In [73]: arr2d[2]  
Out[73]: array([7, 8, 9])
```

Das bedeutet, dass auf einzelne Elemente rekursiv zugegriffen werden kann. Allerdings ist das ein bisschen zu viel Arbeit, Sie können deshalb eine kommaseparierte Liste mit Indizes übergeben, um einzelne Elemente auszuwählen. Diese sind also äquivalent:

```
In [74]: arr2d[0][2]
Out[74]: 3
```

```
In [75]: arr2d[0, 2]
Out[75]: 3
```

Abbildung 4-1 verdeutlicht das Indizieren eines zweidimensionalen Arrays. Ich persönlich finde es hilfreich, mir die Achse 0 als die »Zeilen« und die Achse 1 als die »Spalten« des Arrays vorzustellen.

		Achse 1		
		0	1	2
Achse 2	0	0,0	0,1	0,2
	1	1,0	1,1	1,2
	2	2,0	2,1	2,2

Abbildung 4-1: Indizieren der Elemente in einem NumPy-Array

Falls Sie in einem mehrdimensionalen Array die späteren Indizes weglassen, ist das zurückgelieferte Objekt ein niedrigerdimensionales ndarray, das aus all den Daten aus den höheren Dimensionen besteht. In dem $2 \times 2 \times 3$ -Array `arr3d`:

```
In [76]: arr3d = np.array([[1, 2, 3], [4, 5, 6]], [[7, 8, 9], [10, 11, 12]])
```

```
In [77]: arr3d
Out[77]:
array([[ [ 1,  2,  3],
         [ 4,  5,  6]],
       [[ 7,  8,  9],
        [10, 11, 12]])
```

ist `arr3d[0]` also ein 2×3 -Array:

```
In [78]: arr3d[0]
Out[78]:
array([[1, 2, 3],
       [4, 5, 6]])
```

`arr3d[0]` können sowohl skalare Werte als auch Arrays zugewiesen werden:

```
In [79]: old_values = arr3d[0].copy()
```

```
In [80]: arr3d[0] = 42
```

```
In [81]: arr3d
```

```
Out[81]:
array([[42, 42, 42],
       [42, 42, 42]],
      [[ 7,  8,  9],
       [10, 11, 12]])
```

```
In [82]: arr3d[0] = old_values
```

```
In [83]: arr3d
Out[83]:
array([[ 1,  2,  3],
       [ 4,  5,  6]],
      [[ 7,  8,  9],
       [10, 11, 12]])
```

In gleicher Weise liefert Ihnen `arr3d[1, 0]` alle Werte, deren Indizes mit `(1, 0)` starten, sodass ein 1-dimensionales Array gebildet wird:

```
In [84]: arr3d[1, 0]
Out[84]: array([7, 8, 9])
```

Hätten wir in zwei Schritten indiziert, wären wir ebenfalls auf diesen Ausdruck gekommen:

```
In [85]: x = arr3d[1]
```

```
In [86]: x
Out[86]:
array([[ 7,  8,  9],
       [10, 11, 12]])
```

```
In [87]: x[0]
Out[87]: array([7, 8, 9])
```

Beachten Sie, dass in all den Fällen, in denen Teilbereiche des Arrays ausgewählt wurden, die zurückgelieferten Arrays Views sind.

Mit Slices indizieren

Wie eindimensionale Objekte (z.B. Python-Listen) können auch `ndarrays` mit der bekannten Syntax zerlegt werden:

```
In [88]: arr
Out[88]: array([ 0,  1,  2,  3,  4, 64, 64, 64,  8,  9])
```

```
In [89]: arr[1:6]
Out[89]: array([ 1,  2,  3,  4, 64])
```

Betrachten Sie das bereits bekannte zweidimensionale Array `arr2d`. Das Slicing dieses Arrays ist ein bisschen anders:

```
In [90]: arr2d
Out[90]:
array([[1, 2, 3],
       [4, 5, 6],
       [7, 8, 9]])
```

```
In [91]: arr2d[:2]
Out[91]:
array([[1, 2, 3],
       [4, 5, 6]])
```

Wie Sie sehen, wurde es entlang der Achse 0 zerteilt, der ersten Achse. Ein Slice wählt also einen Bereich von Elementen entlang einer Achse aus. Es hilft vielleicht, wenn man den Ausdruck `arr2d[:2]` als »wähle die ersten beiden Zeilen von `arr2d`« liest.

Genau wie Sie mehrere Indizes übergeben können, können Sie auch mehrere Slices übergeben:

```
In [92]: arr2d[:2, 1:]
Out[92]:
array([[2, 3],
       [5, 6]])
```

Wenn Sie das Slicing so durchführen, bekommen Sie immer Array-Views mit der gleichen Anzahl von Dimensionen. Durch das Mischen von Integer-Indizes und Slices erhalten Sie Slices mit weniger Dimensionen.

Ich kann zum Beispiel die zweite Zeile, aber nur die ersten beiden Spalten auswählen:

```
In [93]: arr2d[1, :2]
Out[93]: array([4, 5])
```

Auf ähnliche Weise kann ich die dritte Spalte, aber nur die ersten beiden Zeilen auswählen:

```
In [94]: arr2d[:2, 2]
Out[94]: array([3, 6])
```

In Abbildung 4-2 finden Sie eine grafische Veranschaulichung des Ganzen. Der allein stehende Doppelpunkt bedeutet übrigens, dass die gesamte Achse genommen wird. Sie können deshalb folgendermaßen ein Slice der höherdimensionalen Achsen erzeugen:

```
In [95]: arr2d[:, :1]
Out[95]:
array([[1],
       [4],
       [7]])
```

Das Zuweisen an einen Slice-Ausdruck weist natürlich die gesamte Auswahl zu:

```
In [96]: arr2d[:2, 1:] = 0

In [97]: arr2d
Out[97]:
array([[1, 0, 0],
       [4, 0, 0],
       [7, 8, 9]])
```

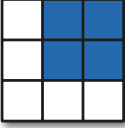
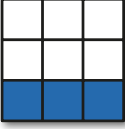
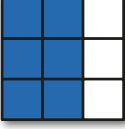
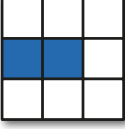
	Ausdruck	Gestalt
	<code>arr[:2, 1:]</code>	(2, 2)
	<code>arr[2]</code> <code>arr[2, :]</code> <code>arr[2:, :]</code>	(3,) (3,) (1, 3)
	<code>arr[:, :2]</code>	(3, 2)
	<code>arr[1, :2]</code> <code>arr[1:2, :2]</code>	(2,) (1, 2)

Abbildung 4-2: Zweidimensionales Slicing von Arrays

Boolesches Indizieren

Betrachten wir ein Beispiel, in dem wir einige Daten in einem Array und ein Array aus mehrfach vorhandenen Namen haben. Ich nutze hier die `randn`-Funktion in `numpy.random`, um einige normalverteilte Daten zu generieren:

```
In [98]: names = np.array(['Bob', 'Joe', 'Will', 'Bob', 'Will', 'Joe', 'Joe'])
```

```
In [99]: data = np.random.randn(7, 4)
```

```
In [100]: names
```

```
Out[100]:
array(['Bob', 'Joe', 'Will', 'Bob', 'Will', 'Joe', 'Joe'],
      dtype='<U4')
```

```
In [101]: data
```

```
Out[101]:
array([[ 0.0929,  0.2817,  0.769 ,  1.2464],
       [ 1.0072, -1.2962,  0.275 ,  0.2289],
       [ 1.3529,  0.8864, -2.0016, -0.3718],
       [ 1.669 , -0.4386, -0.5397,  0.477 ],
       [ 3.2489, -1.0212, -0.5771,  0.1241],
       [ 0.3026,  0.5238,  0.0009,  1.3438],
       [-0.7135, -0.8312, -2.3702, -1.8608]])
```

Nehmen wir an, jeder Name entspricht einer Zeile im Array `data`. Wir wollen alle Zeilen mit dem korrespondierenden Namen 'Bob' auswählen. Genau wie die arithmetischen Operationen sind auch Vergleiche (wie etwa `==`) mit Arrays vektorisiert. Das bedeutet, ein Vergleich von `names` mit dem String 'Bob' ergibt ein boolesches Array:

```
In [102]: names == 'Bob'
Out[102]: array([ True, False, False,  True, False, False, False], dtype=bool)
```

Dieses boolesche Array kann beim Indizieren des Arrays übergeben werden:

```
In [103]: data[names == 'Bob']
Out[103]:
array([[ 0.0929,  0.2817,  0.769 ,  1.2464],
       [ 1.669 , -0.4386, -0.5397,  0.477 ]])
```

Das boolesche Array muss die gleiche Länge haben wie die Array-Achse, die indiziert werden soll. Sie können boolesche Arrays sogar mit Slices und Integer-Werten (oder Sequenzen von Integer-Werten, mehr dazu später) kombinieren.

In diesen Beispielen wähle ich aus den Zeilen aus, in denen `names == 'Bob'` ist, und indiziere außerdem die Spalten:

```
In [104]: data[names == 'Bob', 2:]
Out[104]:
array([[ 0.769 ,  1.2464],
       [-0.5397,  0.477 ]])
```

```
In [105]: data[names == 'Bob', 3]
Out[105]: array([ 1.2464,  0.477 ])
```

Um alles bis auf 'Bob' auszuwählen, können Sie entweder `!=` einsetzen oder die Bedingung mit `~` negieren:

```
In [106]: names != 'Bob'
Out[106]: array([False,  True,  True, False,  True,  True,  True], dtype=bool)
```

```
In [107]: data[~(names == 'Bob')]
Out[107]:
array([[ 1.0072, -1.2962,  0.275 ,  0.2289],
       [ 1.3529,  0.8864, -2.0016, -0.3718],
       [ 3.2489, -1.0212, -0.5771,  0.1241],
       [ 0.3026,  0.5238,  0.0009,  1.3438],
       [-0.7135, -0.8312, -2.3702, -1.8608]])
```

Der Operator `~` ist ganz nützlich, wenn Sie eine allgemeine Bedingung invertieren wollen:

```
In [108]: cond = names == 'Bob'

In [109]: data[~cond]
Out[109]:
array([[ 1.0072, -1.2962,  0.275 ,  0.2289],
       [ 1.3529,  0.8864, -2.0016, -0.3718],
       [ 3.2489, -1.0212, -0.5771,  0.1241],
```



```
[ 0.3026,  0.5238,  0.0009,  1.3438],  
[-0.7135, -0.8312, -2.3702, -1.8608]])
```

Zum Auswählen von zwei der drei Namen nutzen Sie boolesche Operatoren wie & (und) und | (oder):

```
In [110]: mask = (names == 'Bob') | (names == 'Will')  
  
In [111]: mask  
Out[111]: array([ True, False,  True,  True,  True, False, False], dtype=bool)  
  
In [112]: data[mask]  
Out[112]:  
array([[ 0.0929,  0.2817,  0.769 ,  1.2464],  
       [ 1.3529,  0.8864, -2.0016, -0.3718],  
       [ 1.669 , -0.4386, -0.5397,  0.477 ],  
       [ 3.2489, -1.0212, -0.5771,  0.1241]])
```

Die Auswahl von Daten aus einem Array durch boolesche Indizierung erzeugt *immer* eine Kopie der Daten, selbst wenn das zurückgegebene Array unverändert ist.



Die Python-Schlüsselwörter `and` und `or` funktionieren nicht mit booleschen Arrays. Verwenden Sie stattdessen `&` (und) und `|` (oder).

Das Setzen von Werten mit booleschen Arrays ergibt sich aus dem gesunden Menschenverstand. Um alle negativen Werte in `data` auf 0 zu setzen, müssen wir lediglich Folgendes tun:

```
In [113]: data[data < 0] = 0  
  
In [114]: data  
Out[114]:  
array([[ 0.0929,  0.2817,  0.769 ,  1.2464],  
       [ 1.0072,  0.    ,  0.275 ,  0.2289],  
       [ 1.3529,  0.8864,  0.    ,  0.    ],  
       [ 1.669 ,  0.    ,  0.    ,  0.477 ],  
       [ 3.2489,  0.    ,  0.    ,  0.1241],  
       [ 0.3026,  0.5238,  0.0009,  1.3438],  
       [ 0.    ,  0.    ,  0.    ,  0.    ]])
```

Ebenso einfach ist es, ganze Zeilen oder Spalten mithilfe eines eindimensionalen booleschen Arrays zu setzen:

```
In [115]: data[names != 'Joe'] = 7  
  
In [116]: data  
Out[116]:  
array([[ 7.    ,  7.    ,  7.    ,  7.    ],  
       [ 1.0072,  0.    ,  0.275 ,  0.2289],  
       [ 7.    ,  7.    ,  7.    ,  7.    ],  
       [ 7.    ,  7.    ,  7.    ,  7.    ]])
```

```
[ 7.    ,  7.    ,  7.    ,  7.    ],  
 [ 0.3026,  0.5238,  0.0009,  1.3438],  
 [ 0.    ,  0.    ,  0.    ,  0.    ]])
```

Wie Sie später sehen werden, lassen sich diese Arten von Operationen auf zweidimensionalen Daten bequem mit pandas erledigen.

Fancy Indexing

Fancy Indexing (so etwas wie »raffiniertes Indizieren«) ist ein von NumPy übernommener Begriff, der das Indizieren mittels Integer-Arrays beschreibt. Nehmen wir an, wir hätten ein 8×4 -Array:

```
In [117]: arr = np.empty((8, 4))
```

```
In [118]: for i in range(8):  
.....:     arr[i] = i
```

```
In [119]: arr
```

```
Out[119]:
```

```
array([[ 0.,  0.,  0.,  0.],  
       [ 1.,  1.,  1.,  1.],  
       [ 2.,  2.,  2.,  2.],  
       [ 3.,  3.,  3.,  3.],  
       [ 4.,  4.,  4.,  4.],  
       [ 5.,  5.,  5.,  5.],  
       [ 6.,  6.,  6.,  6.],  
       [ 7.,  7.,  7.,  7.]])
```

Um eine Teilmenge der Zeilen in einer bestimmten Reihenfolge auszuwählen, können Sie einfach eine Liste oder ein ndarray mit Integer-Werten übergeben, die die gewünschte Reihenfolge angeben:

```
In [120]: arr[[4, 3, 0, 6]]
```

```
Out[120]:
```

```
array([[ 4.,  4.,  4.,  4.],  
       [ 3.,  3.,  3.,  3.],  
       [ 0.,  0.,  0.,  0.],  
       [ 6.,  6.,  6.,  6.]])
```

Dieser Code hat hoffentlich das gemacht, was Sie erwartet haben! Negative Indizes wählen Zeilen vom Ende her aus:

```
In [121]: arr[[-3, -5, -7]]
```

```
Out[121]:
```

```
array([[ 5.,  5.,  5.,  5.],  
       [ 3.,  3.,  3.,  3.],  
       [ 1.,  1.,  1.,  1.]])
```

Wenn Sie mehrere Index-Arrays übergeben, passiert etwas geringfügig anderes; dabei wird ein eindimensionales Array aus Elementen ausgewählt, die den einzelnen Tupeln der Indizes entsprechen:

```
In [122]: arr = np.arange(32).reshape((8, 4))
```

```
In [123]: arr
Out[123]:
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11],
       [12, 13, 14, 15],
       [16, 17, 18, 19],
       [20, 21, 22, 23],
       [24, 25, 26, 27],
       [28, 29, 30, 31]])
```

```
In [124]: arr[[1, 5, 7, 2], [0, 3, 1, 2]]
Out[124]: array([ 4, 23, 29, 10])
```

Wir schauen uns die `reshape`-Methode in Anhang A genauer an.

Hier wurden die Elemente (1, 0), (5, 3), (7, 1) und (2, 2) ausgewählt. Egal wie viele Dimensionen das Array hat (hier sind es nur 2) – das Ergebnis des Fancy Indexing mit mehreren Integer-Arrays ist immer eindimensional.

Das Verhalten von Fancy Indexing ist in diesem Fall etwas anders, als man es erwarten könnte (mich selbst eingeschlossen), nämlich ein rechteckiger Bereich gebildet aus einer Teilmenge an Zeilen und Spalten der Matrix. Hier ist eine Möglichkeit, das zu erreichen:

```
In [125]: arr[[1, 5, 7, 2]][:, [0, 3, 1, 2]]
Out[125]:
array([[ 4,  7,  5,  6],
       [20, 23, 21, 22],
       [28, 31, 29, 30],
       [ 8, 11,  9, 10]])
```

Denken Sie daran, dass Fancy Indexing im Gegensatz zum Slicing die Daten immer in ein neues Array kopiert.

Arrays transponieren und Achsen tauschen

Transponieren ist eine besondere Art des Umformens, die ebenfalls einen View der zugrunde liegenden Daten zurückliefert, ohne etwas zu kopieren. Arrays besitzen die `transpose`-Methode sowie das besondere `T`-Attribut:

```
In [126]: arr = np.arange(15).reshape((3, 5))
```

```
In [127]: arr
Out[127]:
array([[ 0,  1,  2,  3,  4],
       [ 5,  6,  7,  8,  9],
       [10, 11, 12, 13, 14]])
```

```
In [128]: arr.T
Out[128]:
array([[ 0,  5, 10],
       [ 1,  6, 11],
```

```
[ 2,  7, 12],
 [ 3,  8, 13],
 [ 4,  9, 14]])
```

Bei Matrixberechnungen werden Sie das recht oft machen – zum Beispiel wenn Sie mit der `np.dot`-Funktion das innere Matrixprodukt berechnen:

```
In [129]: arr = np.random.randn(6, 3)
```

```
In [130]: arr
```

```
Out[130]:
```

```
array([[ -0.8608,  0.5601, -1.2659],
       [ 0.1198, -1.0635,  0.3329],
       [-2.3594, -0.1995, -1.542 ],
       [-0.9707, -1.307 ,  0.2863],
       [ 0.378 , -0.7539,  0.3313],
       [ 1.3497,  0.0699,  0.2467]])
```

```
In [131]: np.dot(arr.T, arr)
```

```
Out[131]:
```

```
array([[ 9.2291,  0.9394,  4.948 ],
       [ 0.9394,  3.7662, -1.3622],
       [ 4.948 , -1.3622,  4.3437]])
```

Für höherdimensionale Arrays akzeptiert `transpose` ein Tupel der Achsennummern, die umgestellt werden sollen:

```
In [132]: arr = np.arange(16).reshape((2, 2, 4))
```

```
In [133]: arr
```

```
Out[133]:
```

```
array([[[[ 0,  1,  2,  3],
         [ 4,  5,  6,  7]],
       [[ 8,  9, 10, 11],
        [12, 13, 14, 15]])])
```

```
In [134]: arr.transpose((1, 0, 2))
```

```
Out[134]:
```

```
array([[[[ 0,  1,  2,  3],
         [ 8,  9, 10, 11]],
       [[ 4,  5,  6,  7],
        [12, 13, 14, 15]])])
```

Hier wurden die Achsen so umsortiert, dass die zweite Achse zuerst kommt und die erste Achse als Zweites, während die letzte Achse unverändert blieb.

Das einfache Transponieren mit `.T` ist ein Sonderfall des Vertauschens von Achsen. `ndarray` besitzt die Methode `swapaxes`, die ein Paar aus Achsennummern entgegennimmt und die angegebenen Achsen vertauscht, um die Daten neu anzuordnen:

```
In [135]: arr
```

```
Out[135]:
```

```
array([[[[ 0,  1,  2,  3],
         [ 4,  5,  6,  7]],
```

```
[[ 8,  9, 10, 11],  
 [12, 13, 14, 15]])
```

```
In [136]: arr.swapaxes(1, 2)  
Out[136]:  
array([[[ 0,  4],  
        [ 1,  5],  
        [ 2,  6],  
        [ 3,  7]],  
       [[ 8, 12],  
        [ 9, 13],  
        [10, 14],  
        [11, 15]])])
```

`swapaxes` liefert ebenfalls einen View der Daten zurück, ohne eine Kopie herzustellen.

4.2 Universelle Funktionen: schnelle elementweise Array-Funktionen

Eine universelle Funktion oder *ufunc* ist eine Funktion, die auf den Daten in `ndarrays` elementweise Operationen durchführt. Stellen Sie sie sich als einen schnellen vektorisierten Wrapper für einfache Funktionen vor, der einen oder mehrere skalare Werte nimmt und ein oder mehrere skalare Ergebnisse produziert.

Viele *ufuncs* sind einfache elementweise Transformationen, wie `sqrt` oder `exp`:

```
In [137]: arr = np.arange(10)  
  
In [138]: arr  
Out[138]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])  
  
In [139]: np.sqrt(arr)  
Out[139]:  
array([ 0.         ,  1.         ,  1.41421356,  1.73205081,  2.         ,  2.23606798,  
        2.44948974,  2.64575131,  3.         ])  
  
In [140]: np.exp(arr)  
Out[140]:  
array([  1.         ,  2.71828183,  7.38906809,  20.08553692,  54.59815003,  
        148.4131591 ,  403.42879349, 1096.63315836, 2980.95798704, 8103.08392757])
```

Diese werden als *unäre* *ufuncs* bezeichnet. Andere, wie `add` oder `maximum`, nehmen zwei Arrays entgegen (daher *binäre* *ufuncs*) und liefern ein einzelnes Array als Ergebnis zurück:

```
In [141]: x = np.random.randn(8)  
  
In [142]: y = np.random.randn(8)  
  
In [143]: x
```

```
Out[143]:
array([-0.0119,  1.0048,  1.3272, -0.9193, -1.5491,  0.0222,  0.7584,
        -0.6605])
```

```
In [144]: y
Out[144]:
array([ 0.8626, -0.01  ,  0.05  ,  0.6702,  0.853  , -0.9559, -0.0235,
        -2.3042])
```

```
In [145]: np.maximum(x, y)
Out[145]:
array([ 0.8626,  1.0048,  1.3272,  0.6702,  0.853  ,  0.0222,  0.7584,
        -0.6605])
```

Hier hat `numpy.maximum` das elementweise Maximum der Elemente in `x` und `y` berechnet.

Es ist zwar nicht üblich, aber eine `ufunc` kann auch mehrere Arrays zurückgeben. `modf` ist ein Beispiel dafür, eine vektorisierte Version der eingebauten Python-Funktion `divmod`; es liefert die gebrochenen und ganzzahligen Teile eines Gleitkomma-Arrays zurück:

```
In [146]: arr = np.random.randn(7) * 5

In [147]: arr
Out[147]: array([-3.2623, -6.0915, -6.663 ,  5.3731,  3.6182,  3.45  ,  5.0077])

In [148]: remainder, whole_part = np.modf(arr)

In [149]: remainder
Out[149]: array([-0.2623, -0.0915, -0.663 ,  0.3731,  0.6182,  0.45  ,  0.0077])

In [150]: whole_part
Out[150]: array([-3., -6., -6.,  5.,  3.,  3.,  5.])
```

`Ufuncs` akzeptieren ein optionales `out`-Argument, das ihnen erlaubt, an Ort und Stelle auf Arrays zu operieren:

```
In [151]: arr
Out[151]: array([-3.2623, -6.0915, -6.663 ,  5.3731,  3.6182,  3.45  ,  5.0077])

In [152]: np.sqrt(arr)
Out[152]: array([  nan,    nan,    nan,  2.318 ,  1.9022,  1.8574,  2.2378])

In [153]: np.sqrt(arr, arr)
Out[153]: array([  nan,    nan,    nan,  2.318 ,  1.9022,  1.8574,  2.2378])

In [154]: arr
Out[154]: array([  nan,    nan,    nan,  2.318 ,  1.9022,  1.8574,  2.2378])
```

In den Tabellen [Tabelle 4-3](#) und [Tabelle 4-4](#) finden Sie eine Auflistung der verfügbaren `ufuncs`.

Tabelle 4-3: Unäre ufuncs

Funktion	Beschreibung
abs, fabs	Berechnet elementweise den absoluten Wert für Integer-, Gleitkomma- oder komplexe Werte.
sqrt	Berechnet für jedes Element die Quadratwurzel (äquivalent zu <code>arr ** 0.5</code>).
square	Berechnet das Quadrat jedes Elements (äquivalent zu <code>arr ** 2</code>).
exp	Berechnet den Exponenten e^x jedes Elements.
log, log10, log2, log1p	Natürlicher Logarithmus (Basis e), log Basis 10, log Basis 2 bzw. $\log(1 + x)$.
sign	Berechnet das Vorzeichen jedes Elements: 1 (positiv), 0 (null) oder -1 (negativ).
ceil	Rundet jedes Element auf die kleinste Integer-Zahl auf, die größer oder gleich diesem Element ist.
floor	Rundet jedes Element auf die größte Integer-Zahl ab, die kleiner oder gleich diesem Element ist.
rint	Rundet Elemente auf die nächstgelegene Integer-Zahl, wobei der <code>dtype</code> beibehalten wird.
modf	Gibt die gebrochenen und ganzzahligen Anteile eines Arrays als separates Array zurück.
isnan	Gibt ein boolesches Array zurück, das angibt, ob die einzelnen Werte NaN (Not a Number, keine Zahl) sind.
isfinite, isinf	Gibt ein boolesches Array zurück, das angibt, ob die einzelnen Werte endlich (nicht- <code>inf</code> , nicht-NaN) bzw. unendlich sind.
cos, cosh, sin, sinh, tan, tanh	Reguläre trigonometrische und Hyperbelfunktionen.
arccos, arccosh, arcsin, arcsinh, arctan, arctanh	Inverse trigonometrische Funktionen.
logical_not	Berechnet elementweise den Wahrheitswert von <code>not x</code> (äquivalent zu <code>~arr</code>).

Tabelle 4-4: Binäre universelle Funktionen

Funktion	Beschreibung
add	Addiert korrespondierende Elemente in Arrays.
subtract	Subtrahiert die Elemente im zweiten Array vom ersten Array.
multiply	Multipliziert die Array-Elemente miteinander.
divide, floor_divide	Dividieren oder Floor-Dividieren (Abschneiden des Rests).
power	Potenziiert Elemente im ersten Array entsprechend dem Exponenten im zweiten Array.
maximum, fmax	Elementweises Maximum; <code>fmax</code> ignoriert NaN.
minimum, fmin	Elementweises Minimum; <code>fmin</code> ignoriert NaN.
mod	Elementweiser Modulus (Rest der Division).
copysign	Kopiert die Vorzeichen der Werte im zweiten Argument auf die Werte im ersten Argument.

Tabelle 4-4: Binäre universelle Funktionen (Fortsetzung)

Funktion	Beschreibung
greater, greater_equal, less, less_equal, equal, not_equal	Führt elementweise Vergleiche durch, die ein boolesches Array ergeben (äquivalent zu den Infix-Operatoren >, >=, <, <=, ==, !=).
logical_and, logical_or, logical_xor	Berechnet elementweise die Wahrheitswerte der logischen Operation (äquivalent zu den Infix-Operatoren & , ^).

4.3 Array-orientierte Programmierung mit Arrays

Mit NumPy-Arrays können Sie viele Arten von Datenverarbeitungsaufgaben in Form präziser Array-Ausdrücke schreiben, für die ansonsten Schleifen erforderlich wären. Dieses Ersetzen expliziter Schleifen durch Array-Ausdrücke bezeichnet man gemeinhin als *Vektorisierung*. Im Allgemeinen sind vektorisierte Array-Operationen um ein oder zwei (oder sogar mehr) Größenordnungen schneller als ihre reinen Python-Gegenstücke. Ihre größte Wirkung zeigen sie bei allen Arten von numerischen Berechnungen. Später in Anhang A erläutere ich das *Broadcasting*, eine leistungsstarke Methode für vektorisierte Berechnungen.

Als ein einfaches Beispiel nehmen wir einmal an, dass wir die Funktion $\sqrt{x^2 + y^2}$ über ein normales Gitter von Werten evaluieren wollten. Die Funktion `np.meshgrid` nimmt zwei eindimensionale Arrays entgegen und erzeugt zwei zweidimensionale Matrizen entsprechend den Paaren von (x, y) in den beiden Arrays:

```
In [155]: points = np.arange(-5, 5, 0.01) # 1000 equally spaced points
```

```
In [156]: xs, ys = np.meshgrid(points, points)
```

```
In [157]: ys
```

```
Out[157]:
```

```
array([[ -5.   , -5.   , -5.   , ..., -5.   , -5.   , -5.   ],
       [-4.99, -4.99, -4.99, ..., -4.99, -4.99, -4.99],
       [-4.98, -4.98, -4.98, ..., -4.98, -4.98, -4.98],
       ...,
       [ 4.97,  4.97,  4.97, ...,  4.97,  4.97,  4.97],
       [ 4.98,  4.98,  4.98, ...,  4.98,  4.98,  4.98],
       [ 4.99,  4.99,  4.99, ...,  4.99,  4.99,  4.99]])
```

Zum Auswerten der Funktion schreiben Sie einfach den gleichen Ausdruck auf, den Sie auch bei zwei Punkten schreiben würden:

```
In [158]: z = np.sqrt(xs ** 2 + ys ** 2)
```

```
In [159]: z
```

```
Out[159]:
```

```
array([[ 7.0711,  7.064 ,  7.0569, ...,  7.0499,  7.0569,  7.064 ],
       [ 7.064 ,  7.0569,  7.0499, ...,  7.0428,  7.0499,  7.0569],
       [ 7.0569,  7.0499,  7.0428, ...,  7.0357,  7.0428,  7.0499]])
```



```
...,
 [ 7.0499,  7.0428,  7.0357, ...,  7.0286,  7.0357,  7.0428],
 [ 7.0569,  7.0499,  7.0428, ...,  7.0357,  7.0428,  7.0499],
 [ 7.064 ,  7.0569,  7.0499, ...,  7.0428,  7.0499,  7.0569]])
```

In einer Vorschau auf Kapitel 9 nutze ich matplotlib, um Visualisierungen dieses zweidimensionalen Arrays herzustellen:

```
In [160]: import matplotlib.pyplot as plt
```

```
In [161]: plt.imshow(z, cmap=plt.cm.gray); plt.colorbar()
Out[161]: <matplotlib.colorbar.Colorbar at 0x7f715e3fa630>
```

```
In [162]: plt.title("Image plot of  $\sqrt{x^2 + y^2}$  for a grid of values")
Out[162]: <matplotlib.text.Text at 0x7f715d2de748>
```

In Abbildung 4-3 habe ich die matplotlib-Funktion `imshow` genommen, um aus dem zweidimensionalen Array von Funktionswerten ein Bild zeichnen zu lassen.

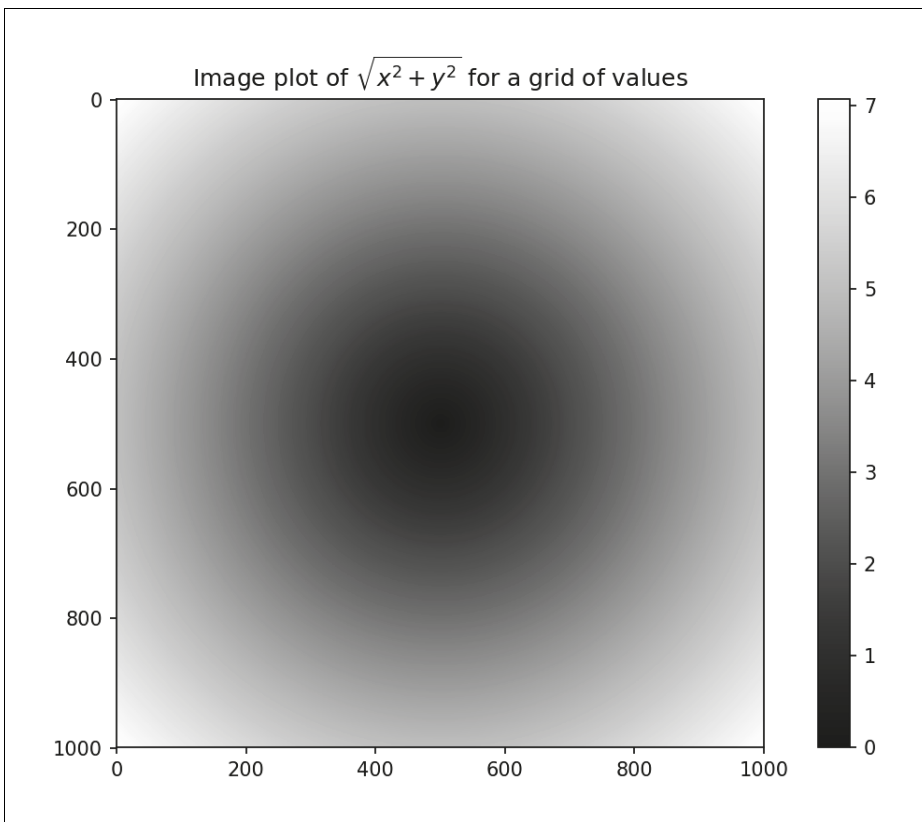


Abbildung 4-3: Plot der Funktionsauswertung am Grid

Bedingte Logik als Array-Operationen ausdrücken

Die Funktion `numpy.where` ist eine vektorisierte Version des ternären Ausdrucks `x if condition else y`. Nehmen wir an, wir hätten ein boolesches Array und zwei Arrays mit Werten:

```
In [165]: xarr = np.array([1.1, 1.2, 1.3, 1.4, 1.5])
```

```
In [166]: yarr = np.array([2.1, 2.2, 2.3, 2.4, 2.5])
```

```
In [167]: cond = np.array([True, False, True, True, False])
```

Nun wollen wir einen Wert aus `xarr` nehmen, wenn der korrespondierende Wert in `cond` `True` ist. Andernfalls soll der Wert aus `yarr` genommen werden. Eine List Comprehension, die das erledigt, könnte so aussehen:

```
In [168]: result = [(x if c else y)
.....:                  for x, y, c in zip(xarr, yarr, cond)]
```

```
In [169]: result
```

```
Out[169]: [1.1000000000000001, 2.2000000000000002, 1.3, 1.3999999999999999, 2.5]
```

Das bringt mehrere Probleme mit sich. Zum einen ist es für große Arrays nicht besonders schnell (weil die ganze Arbeit in interpretiertem Python-Code ausgeführt wird), und zum anderen funktioniert es nicht mit mehrdimensionalen Arrays. Mit `np.where` können Sie dies sehr knapp und prägnant schreiben:

```
In [170]: result = np.where(cond, xarr, yarr)
```

```
In [171]: result
```

```
Out[171]: array([ 1.1,  2.2,  1.3,  1.4,  2.5])
```

Das zweite bzw. dritte Argument zu `np.where` muss kein Array sein, eines oder beide können Skalare sein. Eine typische Anwendung von `where` in der Datenanalyse ist es, ein neues Array aus Werten auf der Basis eines anderen Arrays zu erzeugen. Stellen Sie sich vor, Sie hätten eine Matrix aus zufällig generierten Daten und wollten alle positiven Werte durch 2 und alle negativen Werte durch `-2` ersetzen. Mit `np.where` ist das ganz einfach:

```
In [172]: arr = np.random.randn(4, 4)
```

```
In [173]: arr
```

```
Out[173]:
```

```
array([[ -0.5031,  -0.6223,  -0.9212,  -0.7262],
       [  0.2229,   0.0513,  -1.1577,   0.8167],
       [  0.4336,   1.0107,   1.8249,  -0.9975],
       [  0.8506,  -0.1316,   0.9124,   0.1882]])
```

```
In [174]: arr > 0
```

```
Out[174]:
```

```
array([[False, False, False, False],
       [ True,  True, False,  True],
```

```
[ True,  True,  True, False],  
 [ True, False,  True,  True]], dtype=bool)
```

```
In [175]: np.where(arr > 0, 2, -2)  
Out[175]:  
array([[ -2,  -2,  -2,  -2],  
       [  2,   2,  -2,   2],  
       [  2,   2,   2,  -2],  
       [  2,  -2,   2,   2]])
```

Wenn Sie `np.where` verwenden, können Sie Skalare und Arrays kombinieren. Zum Beispiel kann ich so alle positiven Werte in `arr` durch die Konstante 2 ersetzen:

```
In [176]: np.where(arr > 0, 2, arr) # set only positive values to 2  
Out[176]:  
array([[ -0.5031,  -0.6223,  -0.9212,  -0.7262],  
       [  2.      ,  2.      , -1.1577,  2.      ],  
       [  2.      ,  2.      ,  2.      , -0.9975],  
       [  2.      , -0.1316,  2.      ,  2.      ]])
```

Die an `np.where` übergebenen Arrays können mehr sein als nur gleich große Arrays oder Skalare.

Mathematische und statistische Methoden

Eine Reihe mathematischer Funktionen, die Statistiken über ein ganzes Array oder über die Daten entlang einer Achse berechnen, sind als Methoden der Array-Klasse verfügbar. Aggregationen (oft auch *Reduktionen* genannt) wie die Summe `sum`, der Mittelwert `mean` und die Standardabweichung `std` können entweder als Methoden der Array-Instanz oder auf oberster Ebene als NumPy-Funktion aufgerufen werden.

Hier generiere ich einige normalverteilte Zufallsdaten und berechne ein paar Statistiken:

```
In [177]: arr = np.random.randn(5, 4)  
  
In [178]: arr  
Out[178]:  
array([[ 2.1695, -0.1149,  2.0037,  0.0296],  
       [ 0.7953,  0.1181, -0.7485,  0.585 ],  
       [ 0.1527, -1.5657, -0.5625, -0.0327],  
       [-0.929 , -0.4826, -0.0363,  1.0954],  
       [ 0.9809, -0.5895,  1.5817, -0.5287]])  
  
In [179]: arr.mean()  
Out[179]: 0.19607051119998253  
  
In [180]: np.mean(arr)  
Out[180]: 0.19607051119998253  
  
In [181]: arr.sum()  
Out[181]: 3.9214102239996507
```

Funktionen wie `mean` und `sum` akzeptieren ein optionales `axis`-Argument, sodass die Statistiken entlang der angegebenen Achse berechnet werden. Das Ergebnis ist ein Array mit einer Dimension weniger:

```
In [182]: arr.mean(axis=1)
Out[182]: array([ 1.022 ,  0.1875, -0.502 , -0.0881,  0.3611])
```

```
In [183]: arr.sum(axis=0)
Out[183]: array([ 3.1693, -2.6345,  2.2381,  1.1486])
```

Hier bedeutet `arr.mean(1)`: »Berechne den Mittelwert über die Spalten.«, während `arr.sum(0)` bedeutet: »Berechne die Summe entlang der Zeilen.«

Andere Methoden wie `cumsum` und `cumprod` aggregieren nicht, sondern erzeugen ein Array der Zwischenergebnisse:

```
In [184]: arr = np.array([0, 1, 2, 3, 4, 5, 6, 7])
```

```
In [185]: arr.cumsum()
Out[185]: array([ 0,  1,  3,  6, 10, 15, 21, 28])
```

In mehrdimensionalen Arrays liefern Akkumulationsfunktionen wie `cumsum` ein Array der gleichen Größe zurück, allerdings mit Teilergebnissen entlang der Achsen entsprechend den niedriger dimensionierten Slices:

```
In [186]: arr = np.array([[0, 1, 2], [3, 4, 5], [6, 7, 8]])
```

```
In [187]: arr
Out[187]:
array([[0, 1, 2],
       [3, 4, 5],
       [6, 7, 8]])
```

```
In [188]: arr.cumsum(axis=0)
Out[188]:
array([[ 0,  1,  2],
       [ 3,  5,  7],
       [ 9, 12, 15]])
```

```
In [189]: arr.cumprod(axis=1)
Out[189]:
array([[ 0,  0,  0],
       [ 3, 12, 60],
       [ 6, 42, 336]])
```

In Tabelle 4-5 finden Sie eine vollständige Liste. In späteren Kapiteln werden wir dann Beispiele dieser Methoden in Aktion sehen.

Tabelle 4-5: Grundlegende statistische Array-Methoden

Methode	Beschreibung
<code>sum</code>	Summe aller Elemente in dem Array oder entlang einer Achse; Arrays der Länge null ergeben die Summe 0.
<code>mean</code>	Arithmetischer Mittelwert; Arrays der Länge null haben den Mittelwert NaN.

Tabelle 4-5: Grundlegende statistische Array-Methoden (Fortsetzung)

Method	Beschreibung
std, var	Standardabweichung bzw. Varianz mit optionaler Anpassung der Freiheitsgrade (Standardnenner n).
min, max	Minimum und Maximum.
argmin, argmax	Indizes der minimalen bzw. maximalen Elemente.
cumsum	Kumulative Summe der Elemente, beginnend bei 0.
cumprod	Kumulatives Produkt der Elemente, beginnend bei 1.

Methoden für boolesche Arrays

Boolesche Werte werden in den vorgestellten Methoden als 1 (True) und 0 (False) behandelt. Daher wird `sum` oft als Möglichkeit genutzt, True-Werte in einem booleschen Array zu zählen:

```
In [190]: arr = np.random.randn(100)

In [191]: (arr > 0).sum() # Number of positive values
Out[191]: 42
```

Es gibt zwei weitere Methoden, `any` und `all`, die besonders bei booleschen Arrays nützlich sind. `any` testet, ob ein oder mehrere Werte in einem Array True sind, während `all` prüft, ob alle Werte True sind:

```
In [192]: bools = np.array([False, False, True, False])

In [193]: bools.any()
Out[193]: True

In [194]: bools.all()
Out[194]: False
```

Diese Methoden funktionieren auch mit nicht booleschen Arrays, bei denen Werte, die nicht 0 sind, zu True ausgewertet werden.

Sortieren

Wie Python's eingebauter Listentyp können auch NumPy-Arrays mit der Methode `sort` an Ort und Stelle sortiert werden:

```
In [195]: arr = np.random.randn(6)

In [196]: arr
Out[196]: array([ 0.6095, -0.4938,  1.24  , -0.1357,  1.43  , -0.8469])

In [197]: arr.sort()

In [198]: arr
Out[198]: array([-0.8469, -0.4938, -0.1357,  0.6095,  1.24  ,  1.43  ])
```

Bei einem mehrdimensionalen Array können Sie jeden eindimensionalen Bereich mit Werten an Ort und Stelle entlang einer Achse sortieren, wenn Sie die Achsennummer an `sort` übergeben:

```
In [199]: arr = np.random.randn(5, 3)
```

```
In [200]: arr
```

```
Out[200]:
```

```
array([[ 0.6033,  1.2636, -0.2555],
       [-0.4457,  0.4684, -0.9616],
       [-1.8245,  0.6254,  1.0229],
       [ 1.1074,  0.0909, -0.3501],
       [ 0.218 , -0.8948, -1.7415]])
```

```
In [201]: arr.sort(1)
```

```
In [202]: arr
```

```
Out[202]:
```

```
array([[ -0.2555,  0.6033,  1.2636],
       [-0.9616, -0.4457,  0.4684],
       [-1.8245,  0.6254,  1.0229],
       [-0.3501,  0.0909,  1.1074],
       [-1.7415, -0.8948,  0.218 ]])
```

Die Methode `np.sort` liefert eine sortierte Kopie eines Arrays zurück, statt das Array direkt zu modifizieren. Ein schneller Weg, die Quantile eines Arrays zu berechnen, besteht darin, es zu sortieren und dann die Werte an einem bestimmten Rang auszuwählen:

```
In [203]: large_arr = np.random.randn(1000)
```

```
In [204]: large_arr.sort()
```

```
In [205]: large_arr[int(0.05 * len(large_arr))] # 5% quantile
```

```
Out[205]: -1.5311513550102103
```

Weitere Details über die Sortiermethoden von NumPy und andere komplexere Techniken wie das indirekte Sortieren finden Sie in Anhang A. In pandas gibt es ebenfalls Arten der Datenmanipulation, die mit dem Sortieren zu tun haben (wie etwa das Sortieren einer Datentabelle nach einer oder mehreren Spalten).

Unique und andere Mengenlogik

NumPy besitzt einige elementare Mengenoperationen für eindimensionale `ndarrays`. Häufig benutzt wird `np.unique`, die die sortierten eindeutigen Werte in einem Array zurückliefert:

```
In [206]: names = np.array(['Bob', 'Joe', 'Will', 'Bob', 'Will', 'Joe', 'Joe'])
```

```
In [207]: np.unique(names)
```

```
Out[207]:
```

```
array(['Bob', 'Joe', 'Will'],
      dtype='<U4')
```

```
In [208]: ints = np.array([3, 3, 3, 2, 2, 1, 1, 4, 4])
```

```
In [209]: np.unique(ints)
Out[209]: array([1, 2, 3, 4])
```

Stellen Sie `np.unique` einmal der reinen Python-Alternative gegenüber:

```
In [210]: sorted(set(names))
Out[210]: ['Bob', 'Joe', 'Will']
```

Eine andere Funktion, `np.in1d`, prüft das Vorhandensein der Werte aus einem Array in einem anderen und liefert ein boolesches Array zurück:

```
In [211]: values = np.array([6, 0, 0, 3, 2, 5, 6])
```

```
In [212]: np.in1d(values, [2, 3, 6])
Out[212]: array([ True, False, False,  True,  True, False,  True], dtype=bool)
```

In Tabelle 4-6 finden Sie eine Liste der Mengenfunktionen in NumPy.

Tabelle 4-6: Mengenoperationen bei Arrays

Methoden	Beschreibung
<code>unique(x)</code>	Berechnet die sortierten eindeutigen Elemente in <code>x</code> .
<code>intersect1d(x, y)</code>	Berechnet die sortierten gemeinsamen Elemente in <code>x</code> und <code>y</code> .
<code>union1d(x, y)</code>	Berechnet die sortierte Vereinigung der Elemente.
<code>in1d(x, y)</code>	Berechnet ein boolesches Array, das angibt, ob die einzelnen Elemente aus <code>x</code> in <code>y</code> enthalten sind.
<code>setdiff1d(x, y)</code>	Mengendifferenz, Elemente in <code>x</code> , die nicht in <code>y</code> enthalten sind.
<code>setxor1d(x, y)</code>	Symmetrische Mengendifferenzen, Elemente, die in einem der Arrays, aber nicht in beiden enthalten sind.

4.4 Dateiein- und -ausgabe bei Arrays

NumPy ist in der Lage, Daten sowohl im Text- als auch im Binärformat auf der Festplatte zu speichern und von ihr zu laden. Ich werde in diesem Abschnitt nur über das in NumPy eingebaute Binärformat sprechen, da die meisten Anwender bevorzugt pandas und andere Werkzeuge zum Laden von Text- oder Tabellendaten benutzen werden (siehe Kapitel 6 für weitere Informationen).

`np.save` und `np.load` sind die zwei Arbeitspferde zum effizienten Speichern und Laden von Array-Daten auf und von der Platte. Arrays werden standardmäßig in einem unkomprimierten rohen Binärformat mit der Dateierweiterung `.npy` gespeichert:

```
In [213]: arr = np.arange(10)
```

```
In [214]: np.save('some_array', arr)
```

Sollte der Dateipfad noch nicht auf `.npy` enden, wird die Erweiterung angehängt. Das Array auf der Platte kann dann mit `np.load` geladen werden:

```
In [215]: np.load('some_array.npy')
Out[215]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

Zum Speichern mehrerer Arrays in einem unkomprimierten Archiv nutzen Sie `np.savez` und übergeben die Arrays als Schlüsselwortargumente:

```
In [216]: np.savez('array_archive.npz', a=arr, b=arr)
```

Wenn Sie eine `.npz`-Datei laden, erhalten Sie ein Dictionary-artiges Objekt zurück, das die einzelnen Arrays dann lädt (»lazy«), wenn auf sie zugegriffen wird:

```
In [217]: arch = np.load('array_archive.npz')
```

```
In [218]: arch['b']
Out[218]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
```

Lassen sich Ihre Daten gut komprimieren, könnten Sie stattdessen `numpy.savez_compressed` benutzen:

```
In [219]: np.savez_compressed('arrays_compressed.npz', a=arr, b=arr)
```

4.5 Lineare Algebra

Lineare Algebra ist, genau wie Matrixmultiplikation, Dekompositionen, Determinanten und andere Mathematik mit quadratischen Matrizen, ein wichtiger Bestandteil jeder Array-Bibliothek. Im Gegensatz zu manchen Sprachen wie etwa MATLAB ist das Multiplizieren zweier zweidimensionaler Arrays mit `*` ein elementweises Produkt und kein Skalarprodukt von Matrizen. Dafür gibt es die Funktion `dot` – sowohl als Array-Methode als auch als Funktion im `numpy`-Namensraum – für die Matrixmultiplikation:

```
In [223]: x = np.array([[1., 2., 3.], [4., 5., 6.]])
```

```
In [224]: y = np.array([[6., 23.], [-1, 7], [8, 9]])
```

```
In [225]: x
Out[225]:
array([[ 1.,  2.,  3.],
       [ 4.,  5.,  6.]])
```

```
In [226]: y
Out[226]:
array([[ 6., 23.],
       [-1.,  7.],
       [ 8.,  9.]])
```

```
In [227]: x.dot(y)
Out[227]:
array([[ 28.,  64.],
       [ 67., 181.]])
```

`x.dot(y)` ist äquivalent zu `np.dot(x, y)`:


```
In [228]: np.dot(x, y)
Out[228]:
array([[ 28.,  64.],
       [ 67., 181.]])
```

Ein Matrixprodukt zwischen einem zweidimensionalen Array und einem eindimensionalen Array passender Größe resultiert in einem eindimensionalen Array:

```
In [229]: np.dot(x, np.ones(3))
Out[229]: array([ 6., 15.])
```

Das @-Symbol (seit Python 3.5) agiert auch als Infix-Operator für Matrixmultiplikationen:

```
In [230]: x @ np.ones(3)
Out[230]: array([ 6., 15.])
```

`numpy.linalg` besitzt einen Standardsatz an Funktionen zu Matrixzerlegung, Invertierungen, Determinanten und anderen Dingen. Die Funktionen sind unter der Oberfläche mit den gleichen Standardbibliotheken für lineare Algebra implementiert, die auch in anderen Sprachen wie MATLAB und R zum Einsatz kommen, also zum Beispiel BLAS, LAPACK oder vielleicht auch (je nach Ihrer NumPy-Version) der proprietären Intel-MKL (Math Kernel Library):

```
In [231]: from numpy.linalg import inv, qr
```

```
In [232]: X = np.random.randn(5, 5)
```

```
In [233]: mat = X.T.dot(X)
```

```
In [234]: inv(mat)
Out[234]:
array([[ 933.1189,  871.8258, -1417.6902, -1460.4005, 1782.1391],
       [ 871.8258,  815.3929, -1325.9965, -1365.9242, 1666.9347],
       [-1417.6902, -1325.9965, 2158.4424, 2222.0191, -2711.6822],
       [-1460.4005, -1365.9242, 2222.0191, 2289.0575, -2793.422 ],
       [1782.1391, 1666.9347, -2711.6822, -2793.422 , 3409.5128]])
```

```
In [235]: mat.dot(inv(mat))
Out[235]:
array([[ 1.,  0., -0., -0., -0.],
       [-0.,  1.,  0.,  0.,  0.],
       [ 0.,  0.,  1.,  0.,  0.],
       [-0.,  0.,  0.,  1., -0.],
       [-0.,  0.,  0.,  0.,  1.]])
```

```
In [236]: q, r = qr(mat)
```

```
In [237]: r
Out[237]:
array([[ -1.6914,  4.38 ,  0.1757,  0.4075, -0.7838],
       [ 0. , -2.6436,  0.1939, -3.072 , -1.0702],
       [ 0. ,  0. , -0.8138,  1.5414,  0.6155],
       [ 0. ,  0. ,  0. , -2.6445, -2.1669],
       [ 0. ,  0. ,  0. ,  0. ,  0.0002]])
```

Der Ausdruck $X.T.dot(X)$ berechnet das Skalarprodukt von X mit seiner transponierten $X.T$.

In Tabelle 4-7 finden Sie eine Liste der gebräuchlichsten Funktionen für die lineare Algebra.

Tabelle 4-7: Häufig verwendete `numpy.linalg`-Funktionen

Funktion	Beschreibung
<code>diag</code>	Gibt die diagonalen (oder gegendiagonalen) Elemente einer quadratischen Matrix als eindimensionales Array zurück oder konvertiert ein eindimensionales Array in eine quadratische Matrix mit Nullen auf der Gegendiagonalen.
<code>dot</code>	Matrixmultiplikation.
<code>trace</code>	Berechnet die Summe der diagonalen Elemente.
<code>det</code>	Berechnet die Matrixdeterminante.
<code>eig</code>	Berechnet die Eigenwerte und Eigenvektoren einer quadratischen Matrix.
<code>inv</code>	Berechnet die Inverse einer quadratischen Matrix.
<code>pinv</code>	Berechnet die Moore-Penrose-Pseudoinverse einer Matrix.
<code>qr</code>	Berechnet die QR-Zerlegung.
<code>svd</code>	Berechnet die Singulärwertzerlegung (Singular Value Decomposition, SVD).
<code>solve</code>	Löst das lineare System $Ax = b$ für x , wobei A eine quadratische Matrix ist.
<code>lstsq</code>	Berechnet die Lösung der kleinsten Quadrate für $Ax = b$.

4.6 Erzeugen von Pseudozufallszahlen

Das Modul `numpy.random` ergänzt das eingebaute Python-Modul `random` um Funktionen, die effizient ganze Arrays mit Zufallswerten aus vielen Arten von Wahrscheinlichkeitsverteilungen erzeugen. So erhalten Sie zum Beispiel mit `normal` ein 4×4 -Array mit Zufallswerten aus der Standardnormalverteilung:

```
In [238]: samples = np.random.normal(size=(4, 4))
```

```
In [239]: samples
```

```
Out[239]:
```

```
array([[ 0.5732,  0.1933,  0.4429,  1.2796],  
       [ 0.575 ,  0.4339, -0.7658, -1.237 ],  
       [-0.5367,  1.8545, -0.92  , -0.1082],  
       [ 0.1525,  0.9435, -1.0953, -0.144 ]])
```

Pythons eingebautes `random`-Modul generiert im Gegensatz dazu immer nur eine Stichprobe auf einmal. Wie Ihnen dieser Benchmark zeigt, ist `numpy.random` beim Generieren sehr großer Stichproben deutlich mehr als eine Größenordnung schneller:

```
In [240]: from random import normalvariate
```

```
In [241]: N = 1000000
```

```
In [242]: %timeit samples = [normalvariate(0, 1) for _ in range(N)]
1.77 s +- 126 ms per loop (mean +- std. dev. of 7 runs, 1 loop each)
```

```
In [243]: %timeit np.random.normal(size=N)
61.7 ms +- 1.32 ms per loop (mean +- std. dev. of 7 runs, 10 loops each)
```

Wir bezeichnen diese Werte als *Pseudozufallszahlen*, weil sie durch einen Algorithmus mit deterministischem Verhalten generiert werden, der auf dem *Startwert* (Seed) des Zufallszahlengenerators beruht. Sie können den Startwert für die Zufallszahlengenerierung mit NumPy mit `np.random.seed` ändern:

```
In [244]: np.random.seed(1234)
```

Die Funktionen zur Datengenerierung in `numpy.random` nutzen einen globalen Zufallsstartwert. Um einen globalen Zustand zu vermeiden, können Sie `numpy.random.RandomState` einsetzen. Damit legen Sie einen Zufallszahlengenerator an, der von anderen isoliert ist:

```
In [245]: rng = np.random.RandomState(1234)
```

```
In [246]: rng.randn(10)
```

```
Out[246]:
array([ 0.4714, -1.191 ,  1.4327, -0.3127, -0.7206,  0.8872,  0.8596,
        -0.6365,  0.0157, -2.2427])
```

In Tabelle 4-8 finden Sie eine unvollständige Liste der in `numpy.random` verfügbaren Funktionen. Im nächsten Abschnitt werde ich Ihnen an einigen Beispielen zeigen, wie Sie mit diesen Funktionen große Arrays mit Stichproben generieren können.

Tabelle 4-8: Unvollständige Liste der `numpy.random`-Funktionen

Funktion	Beschreibung
<code>seed</code>	Setzt den Startwert des Zufallszahlengenerators.
<code>permutation</code>	Liefert eine zufällige Permutation einer Sequenz oder einen permutierten Bereich zurück.
<code>shuffle</code>	Vermischt eine Sequenz an Ort und Stelle.
<code>rand</code>	Zieht Stichproben aus einer Gleichverteilung.
<code>randint</code>	Zieht zufällige Ganzzahlen aus einem gegebenen Intervall.
<code>randn</code>	Zieht Stichproben aus einer Normalverteilung mit dem Mittelwert 0 und der Standardabweichung 1 (Interface wie bei MATLAB).
<code>binomial</code>	Zieht Stichproben aus einer Binomialverteilung.
<code>normal</code>	Zieht Stichproben aus einer Normalverteilung (einer gaussischen Verteilung).
<code>beta</code>	Zieht Stichproben aus einer Beta-Verteilung.
<code>chisquare</code>	Zieht Stichproben aus einer Chi-Quadrat-Verteilung.
<code>gamma</code>	Zieht Stichproben aus einer Gamma-Verteilung.
<code>uniform</code>	Zieht Stichproben aus einer [0, 1)-Gleichverteilung.

4.7 Beispiel: Random Walks

Die Simulation von Random Walks (https://en.wikipedia.org/wiki/Random_walk) (Irrfahrten) zeigt anschaulich die Anwendung von Array-Operationen. Schauen wir uns zuerst einen einfachen Random Walk an, der bei 0 beginnt und bei dem die Schritte 1 und -1 mit gleicher Wahrscheinlichkeit auftreten.

Und so würde man rein in Python einen Random Walk mit 1.000 Schritten implementieren: mit dem eingebauten `random`-Modul:

```
In [247]: import random
.....: position = 0
.....: walk = [position]
.....: steps = 1000
.....: for i in range(steps):
.....:     step = 1 if random.randint(0, 1) else -1
.....:     position += step
.....:     walk.append(position)
.....:
```

In Abbildung 4-4 sehen Sie einen Beispielpplot der ersten 100 Werte eines dieser Random Walks:

```
In [249]: plt.plot(walk[:100])
```

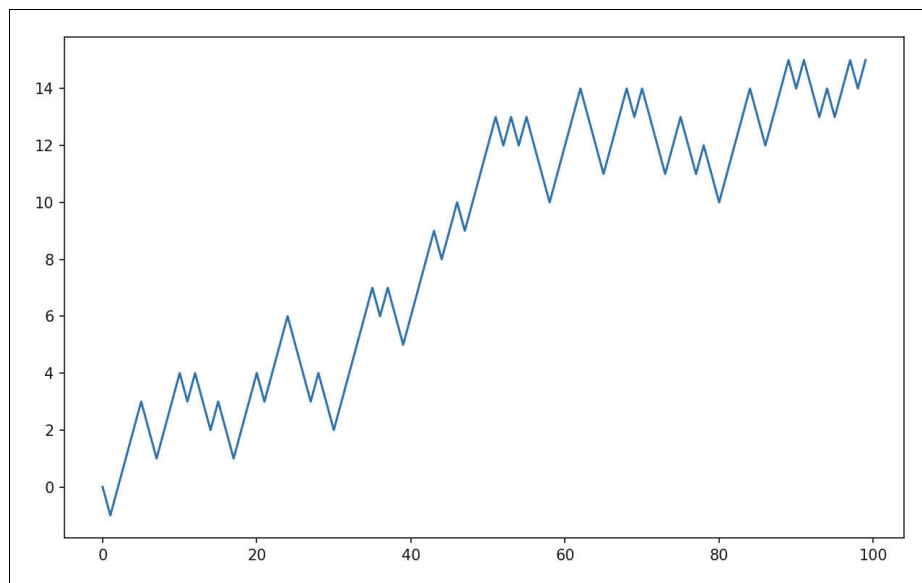


Abbildung 4-4: Eine einfache Irrfahrt

Sie werden vielleicht bemerken, dass `walk` einfach die kumulative Summe der zufälligen Schritte ist und als Array-Ausdruck ausgewertet werden könnte. Daher benutze ich das Modul `np.random`, um auf einmal 1.000 Münzwürfe zu ziehen, setze diese auf 1 und -1 und berechne die kumulative Summe:

```
In [251]: nsteps = 1000
In [252]: draws = np.random.randint(0, 2, size=nsteps)
In [253]: steps = np.where(draws > 0, 1, -1)
In [254]: walk = steps.cumsum()
```

Daraus können wir dann Statistiken wie das Minimum und das Maximum entlang der Bahnkurve ermitteln:

```
In [255]: walk.min()
Out[255]: -3

In [256]: walk.max()
Out[256]: 31
```

Eine kompliziertere Statistik ist die *erste Durchgangszeit*, der Schritt, an dem der Random Walk einen bestimmten Wert erreicht. Hier wollen wir vielleicht wissen, wie lang der Random Walk brauchte, um in jeder Richtung wenigstens 10 Schritte vom Ursprung 0 weg zu sein. `np.abs(walk) >= 10` liefert uns ein boolesches Array, das anzeigt, wo der Walk die 10 erreicht oder überschritten hat, allerdings wollen wir den Index der *ersten* 10 oder -10 . Das können wir mit `argmax` berechnen, das den ersten Index des Maximumwerts im booleschen Array zurückgibt (True ist das Maximum):

```
In [257]: (np.abs(walk) >= 10).argmax()
Out[257]: 37
```

Beachten Sie, dass `argmax` hier nicht immer effizient ist, weil es stets einen vollständigen Scan des Arrays durchführt. In diesem speziellen Fall wissen wir, dass es das Maximum ist, wenn wir einmal ein True festgestellt haben.

Viele Random Walks auf einmal simulieren

Wäre es unser Ziel, viele Random Walks, zum Beispiel 5.000, zu simulieren, könnten Sie mit nur geringen Änderungen an dem gezeigten Code alle Random Walks gleichzeitig generieren. Beim Übergeben eines 2-Tupels generieren die `numpy.random`-Funktionen ein zweidimensionales Array aus Ziehungen, und wir können die kumulative Summe über die Zeilen berechnen, um alle 5.000 Random Walks auf einen Schlag auszurechnen:

```
In [258]: nwalks = 5000
In [259]: nsteps = 1000
In [260]: draws = np.random.randint(0, 2, size=(nwalks, nsteps)) # 0 or 1
In [261]: steps = np.where(draws > 0, 1, -1)
In [262]: walks = steps.cumsum(1)
```

```
In [263]: walks
Out[263]:
array([[ 1,  0,  1, ...,  8,  7,  8],
       [ 1,  0, -1, ..., 34, 33, 32],
       [ 1,  0, -1, ...,  4,  5,  4],
       ...,
       [ 1,  2,  1, ..., 24, 25, 26],
       [ 1,  2,  3, ..., 14, 13, 14],
       [-1, -2, -3, ..., -24, -23, -22]])
```

Jetzt berechnen wir die maximalen und minimalen Werte über alle Walks:

```
In [264]: walks.max()
Out[264]: 138
```

```
In [265]: walks.min()
Out[265]: -133
```

Ermitteln wir nun aus all diesen Walks die minimale Durchgangszeit bis 30 oder -30. Das ist ein bisschen schwieriger, weil nicht alle 5.000 bis 30 gekommen sind. Wir können das mit der `any`-Methode prüfen:

```
In [266]: hits30 = (np.abs(walks) >= 30).any(1)
```

```
In [267]: hits30
Out[267]: array([False,  True, False, ..., False,  True, False], dtype=bool)
```

```
In [268]: hits30.sum() # Number that hit 30 or -30
Out[268]: 3410
```

Mit diesem booleschen Array können wir nun die Zeilen der Walks auswählen, die tatsächlich die absolute 30 erreicht haben. Dann rufen wir `argmax` über die Achse 1 auf, damit die Durchgangszeiten festgestellt werden:

```
In [269]: crossing_times = (np.abs(walks[hits30]) >= 30).argmax(1)
```

Schließlich berechnen wir die minimale durchschnittliche Durchgangszeit:

```
In [270]: crossing_times.mean()
Out[270]: 498.88973607038122
```

Experimentieren Sie selbst mit anderen Verteilungen als nur den gleichverteilten Münzwürfen für die Schritte. Sie brauchen lediglich eine andere Funktion für die Zufallszahlengenerierung, wie etwa `normal`, um normalverteilte Schritte mit einstellbarem Mittelwert und Standardabweichung zu generieren:

```
In [271]: steps = np.random.normal(loc=0, scale=0.25,
.....:                               size=(nwalks, nsteps))
```

4.8 Schlussbemerkung

Der Rest des Buchs konzentriert sich zwar vorrangig auf die Datenverarbeitung mit `pandas`, doch wir werden auch weiter in einem ähnlichen Array-basierten Stil arbeiten. In Anhang A steigen wir tiefer in die `NumPy`-Funktionalität ein, sodass Sie Ihre Fähigkeiten beim Rechnen mit Arrays weiterentwickeln können.