

NumPy avancé

Konrad HINSEN

Centre de Biophysique Moléculaire, Orléans
et
Synchrotron SOLEIL, Saint Aubin

3 décembre 2013

Packaging avec NumPy

Vérifier si NumPy est installé (dans setup.py)

```
1 try:
2     import numpy
3 except ImportError:
4     numpy = None
```

Tester la version

```
1 if [int(s) for s in numpy.__version__.split('.')] \
2     >= [1, 5, 0]:
3     pass
```

Indiquer une dépendance de NumPy (setuptools/distribute)

```
1 setup(...
2     install_requires = ['numpy>1.1'],
3     ...)
```

Où se trouve numpy/arrayobject.h ?

```
1 import numpy.distutils.misc_util
2 numpy_include_dirs = numpy.distutils.misc_util.\
3     get_numpy_include_dirs()
```

Pour les modules C/Cython qui utilisent NumPy:

```
1 Extension (...
2     include_dirs=numpy_include_dirs,
3     ...)
```

Une *vue* est un tableau qui fait référence aux éléments d'un autre tableau au lieu d'avoir son propre espace mémoire. L'indexation crée des vues, tout comme `reshape` et quelques autres méthodes.

Une *vue* est un tableau qui fait référence aux éléments d'un autre tableau au lieu d'avoir son propre espace mémoire. L'indexation crée des vues, tout comme `reshape` et quelques autres méthodes.

Une vue a un attribut `base` dont la valeur est le tableau initial (qui peut lui-même être une vue). Pour un tableau qui n'est pas une vue, `base` est `None`.

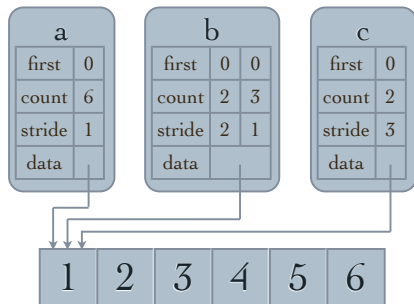
Une *vue* est un tableau qui fait référence aux éléments d'un autre tableau au lieu d'avoir son propre espace mémoire. L'indexation crée des vues, tout comme `reshape` et quelques autres méthodes.

Une vue a un attribut `base` dont la valeur est le tableau initial (qui peut lui-même être une vue). Pour un tableau qui n'est pas une vue, `base` est `None`.

Attention: Les éléments d'un tableau peuvent être modifiés en passant par n'importe quelle vue directe ou indirecte. Il n'est pas possible de vérifier s'il existe des vues pour un tableau donné.

La structure des tableaux dans la mémoire

```
1 >>> import numpy as np
2 >>> a = np.arange(1,7)
3 >>> a
4 array([1, 2, 3, 4, 5, 6])
5 >>> b = a.reshape((3, 2))
6 >>> b
7 array([[1, 2],
8        [3, 4],
9        [5, 6]])
10 >>> c = a[::3]
11 >>> c
12 array([1, 4])
```



```
1 >>> a[0] = 42
2 >>> c
3 array([42, 4])
```



Tant qu'au moins une vue existe, un tableau ne peut pas être détruit.
Après

```
gross = N.arange(100000000)
petit = gross[:5]
del gross
```

le gros tableau reste en mémoire jusqu'à la disparation du petit.
Mieux vaut faire une copie:

```
gross = N.arange(100000000)
petit = gross[:5].copy()
del gross
```


Analyse d'un tableau

Tableau a

```
1 >>> a.dtype
2 dtype('int64')
3 >>> a.itemsize
4 8
5 >>> a.strides
6 (8,)
7 >>> a.__array_interface__ \
8     ['data'][0]
9 4314532672
10 >>> a.flags
11 C_CONTIGUOUS : True
12 F_CONTIGUOUS : True
13 OWNDATA : True
14 WRITEABLE : True
15 ALIGNED : True
16 UPDATEIFCOPY : False
17 >>> a.base is None
18 True
```

Tableau b

```
1 >>> b.dtype
2 dtype('int64')
3 >>> b.itemsize
4 8
5 >>> b.strides
6 (16, 8)
7 >>> b.__array_interface__ \
8     ['data'][0]
9 4314532672
10 >>> b.flags
11 C_CONTIGUOUS : True
12 F_CONTIGUOUS : False
13 OWNDATA : False
14 WRITEABLE : True
15 ALIGNED : True
16 UPDATEIFCOPY : False
17 >>> b.base is a
18 True
```

Tableau c

```
1 >>> c.dtype
2 dtype('int64')
3 >>> c.itemsize
4 8
5 >>> c.strides
6 (24,)
7 >>> c.__array_interface__ \
8     ['data'][0]
9 4314532672
10 >>> c.flags
11 C_CONTIGUOUS : False
12 F_CONTIGUOUS : False
13 OWNDATA : False
14 WRITEABLE : True
15 ALIGNED : True
16 UPDATEIFCOPY : False
17 >>> c.base is a
18 True
```

Tableaux C et Fortran

C

```
1 >>> C = N.array([[1, 2],
2                   [3, 4]],
3                   order='C')
4 >>> C
5 array([[1, 2],
6        [3, 4]])
7 >>> C.strides
8 (16, 8)
9 >>> C.flags
10 C_CONTIGUOUS : True
11 F_CONTIGUOUS : False
12 OWNDATA : True
13 WRITEABLE : True
14 ALIGNED : True
15 UPDATEIFCOPY : False
```

Fortran

```
1 >>> F = N.array([[1, 2],
2                   [3, 4]],
3                   order='F')
4 >>> F
5 array([[1, 2],
6        [3, 4]])
7 >>> F.strides
8 (8, 16)
9 >>> F.flags
10 C_CONTIGUOUS : False
11 F_CONTIGUOUS : True
12 OWNDATA : True
13 WRITEABLE : True
14 ALIGNED : True
15 UPDATEIFCOPY : False
```

Utilisation de la mémoire

Un des inconvénients de NumPy est l'allocation d'un grand nombre de tableaux temporaires. Une expression simple comme $c = 2*a + 3*b$ est en fait traitée comme

```
1 tmp1 = 2*a
2 tmp2 = 3*b
3 c = tmp1 + tmp2
```

La création de trois tableaux, dont deux sont presque immédiatement supprimés, est coûteuse en temps, surtout pour des petits tableaux, et coûteuse en espace mémoire, surtout pour les grands tableaux.

Une autre inconvénient est la mauvaise utilisation du cache du processeur pour ce genre d'opération.

Arithmétique en place

```
1 >>> a = np.arange(10)
2 >>> a
3 array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
4 >>> a *= 2
5 >>> a
6 array([ 0,  2,  4,  6,  8, 10, 12, 14, 16, 18])
7 >>> np.multiply(a, 2, a)
8 array([ 0,  4,  8, 12, 16, 20, 24, 28, 32, 36])
9 >>> a
10 array([ 0,  4,  8, 12, 16, 20, 24, 28, 32, 36])
11 >>> b = np.zeros(a.shape, a.dtype)
12 >>> np.multiply(a, 2, b)
13 array([ 0,  8, 16, 24, 32, 40, 48, 56, 64, 72])
14 >>> a
15 array([ 0,  4,  8, 12, 16, 20, 24, 28, 32, 36])
16 >>> b
17 array([ 0,  8, 16, 24, 32, 40, 48, 56, 64, 72])
```

La bibliothèque numexpr évalue des expressions contenant des tableaux de façon optimisée en les compilant en langage C:

- élimination des tableaux intermédiaires
- accès aux tableaux optimisé pour utilise le cache du processeur
- utilisation de processeurs multiples

Exemple

```
1 import numpy as np
2 import numexpr
3 >>> a = 3.*np.arange(10000)-2.
4 >>> b = 0.1*np.arange(10000)**2
5 >>> numexpr.evaluate("3*a**2-5.*b")
6 array([ 1.20000000e+01,  2.50000000e+00,
7         4.60000000e+01, ...,
8         2.64805036e+09,  2.64858019e+09,  2.64911007e+09])
```

Le type des éléments d'un tableau (dtype) est normalement un type "simple" (`numpy.int`, `numpy.float32`, `numpy.complex64`, ...). Mais le dtype peut aussi être:

- un tuple
- une structure avec des champs nommés
- un tableau

Applications:

- stockage compact de données hétérogènes
- implémenter des modèles de données plus structurés
- interfacer avec des structs en C
- lire/écrire des fichiers HDF5

Éléments tuples

```
1 >>> x = np.zeros((2,), dtype=('i4, f4, a10'))
2 >>> x[:] = [(1, 2., 'Hello'), (2, 3., "World")]
3 >>> x
4 array([(1, 2.0, 'Hello'), (2, 3.0, 'World')],
5        dtype=[('f0', '>i4'), ('f1', '>f4'), ('f2', '|S10')])
6 >>> x[1]
7 (2, 3., "World")
```

Éléments champs nommés

```
1 >>> x = np.zeros(3, dtype={'names':['col1', 'col2'],
2                               'formats':['i4', 'f4']})
3 >>> x
4 array([(0, 0.0), (0, 0.0), (0, 0.0)],
5        dtype=[('col1', '>i4'), ('col2', '>f4')])
6 >>> x[0]
7 (0, 0.0)
8 >>> x['col1']
9 array([0, 0, 0], dtype=int32)
```


Éléments tableaux

```
1 >>> x = np.zeros(3, dtype='3int8, float32, (2,3)float64')
2 >>> x['f0'] = array([1, 2, 3])
3 >>> x.shape
4 (3,)
5 >>> x['f0']
6 array([[1, 2, 3],
7        [1, 2, 3],
8        [1, 2, 3]], dtype=int8)
9 >>> x['f0'].shape
10 (3, 3)
11 >>> x[0]
12 ([1, 2, 3], 0.0, [[0.0, 0.0, 0.0], [0.0, 0.0, 0.0]])
13 >>> x['f2'].shape
14 (3, 2, 3)
15 >>> x[1]['f2'].shape
16 (2, 3)
```