

Numpy and Pandas

Data Science Tools 1

Fall 2021



UNIVERSITY *of*
DENVER

- Numpy and Pandas
- Review Wk4 Homework
- Please wait in the zoom waiting room in the office hour
- Mid-term is next week



- Numpy
- Pandas

- Library for vectorized computation
 - What does this mean?
 - Why is it so efficient?

| Vectorized - add two arrays | Non-vectorized - add two arrays |
|--|--|
| <pre>def add_two_vec(a, b): a = np.array(a) b = np.array(b) return a + b</pre> | <pre>def add_two_slow(a, b): c = [] for i in range(len(a)): c.append(a[i] + b[i]) return c</pre> |

extensively in Pandas, SciPy, Matplotlib, scikit-learn, scikit-image
and most other data science and scientific Python packages

In order to understand this, we need to point out the stark difference between ndarray and list -

What is the difference?

In order to understand this, we need to point out the stark difference between ndarray and list -

Lists are non-homogenous and ndarrays are homogenous

The homogenous nature of ndarrays provides very efficient memory allocation

Instantiation

| | |
|---------------------------------------|--|
| <code>obj = np.array(list)</code> | |
| <code>np.arange()</code> | |
| <code>np.ones(number or shape)</code> | |
| <code>np.eye()</code> | |
| <code>np.random.randn()</code> | |

Properties

| | |
|------------------------|--|
| <code>obj.shape</code> | |
| <code>obj.dtype</code> | |

Slicing and boolean indexing

| | |
|----------------------------------|--|
| <code>obj[0:2,1:3]</code> | |
| <code>obj[[obj > 0.5]]</code> | |
| | |
| | |
| | |

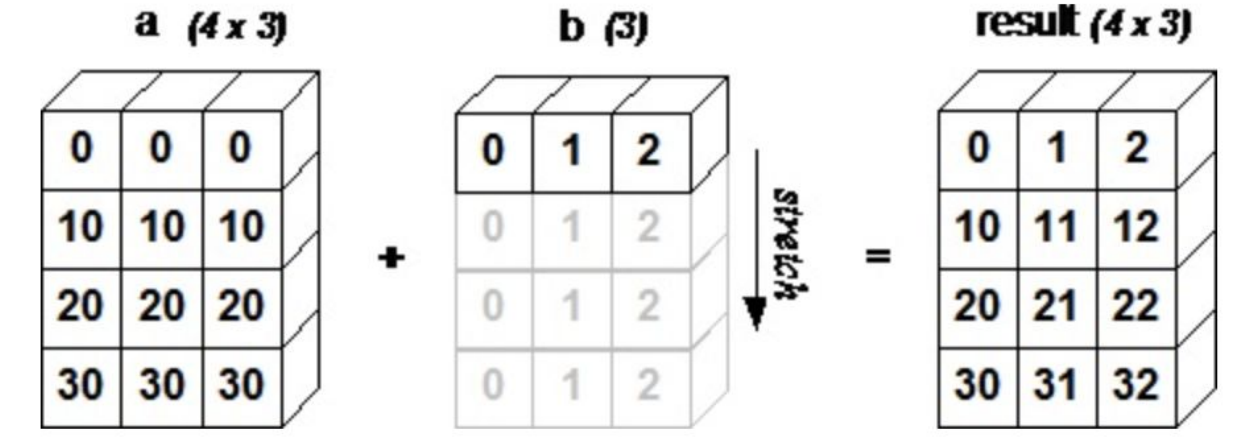
Mathematical Operation

| | |
|---|--|
| <code>obj1 * obj2</code> | element-wise multiplication |
| <code>obj ** obj2</code> | element-wise power |
| <code>np.dot(obj1, obj2)</code> | matrix multiplication |
| <code>obj.T</code> | transpose |
| <code>np.where(cond, obj1, obj2)</code> | |
| <code>np.where(obj>0, 1, -1)</code> | |
| <code>np.mean(obj, axis=0)</code> | 0 = row, 1 = col => <code>[[1,2], [3,4]] = [2, 3]</code> |
| <code>np.argmax</code> | max index |
| <code>np.argmin</code> | min index |
| <code>np.linalg.inv(obj)</code> | inverse |
| <code>np.linalg.svd(obj)</code> | singular value decomposition |

Addresses a limitation of arithmetic operation of two arrays

Broadcasting allows arithmetic operation of different shapes and sizes

Broadcasting solves the problem of arithmetic between arrays of differing shapes by in effect replicating the smaller array along the last mismatched dimension



Numpy - what is broadcasting?

```
import numpy as np  
a = np.array([1, 2, 3])  
b = 2  
c = a + b
```

What is c?

[3, 4,5]

```
import numpy as np  
A = np.array([[1, 2, 3], [1, 2, 3]]) # (2,3)  
b = 2  
C = A + b
```

What is c?

```
[[3, 4, 5], [3, 4, 5]]
```

```
import numpy as np  
A = np.array([[1, 2, 3], [1, 2, 3]]) # (2,3)  
b = np.array([1, 2, 3])  
C = A + b
```

What is c?

```
[[2, 4, 6], [2,4,6]]
```

```
import numpy as np  
A = np.array([[1, 2, 3], [1, 2, 3]])  
b = np.array([1, 2])  
C = A + b
```

What is c?

Error out


```
import numpy as np
A = np.array([[1, 2, 3], [1, 2, 3]])
b = np.array([1, 2]) # b.shape => (2,)
C = A + b
```

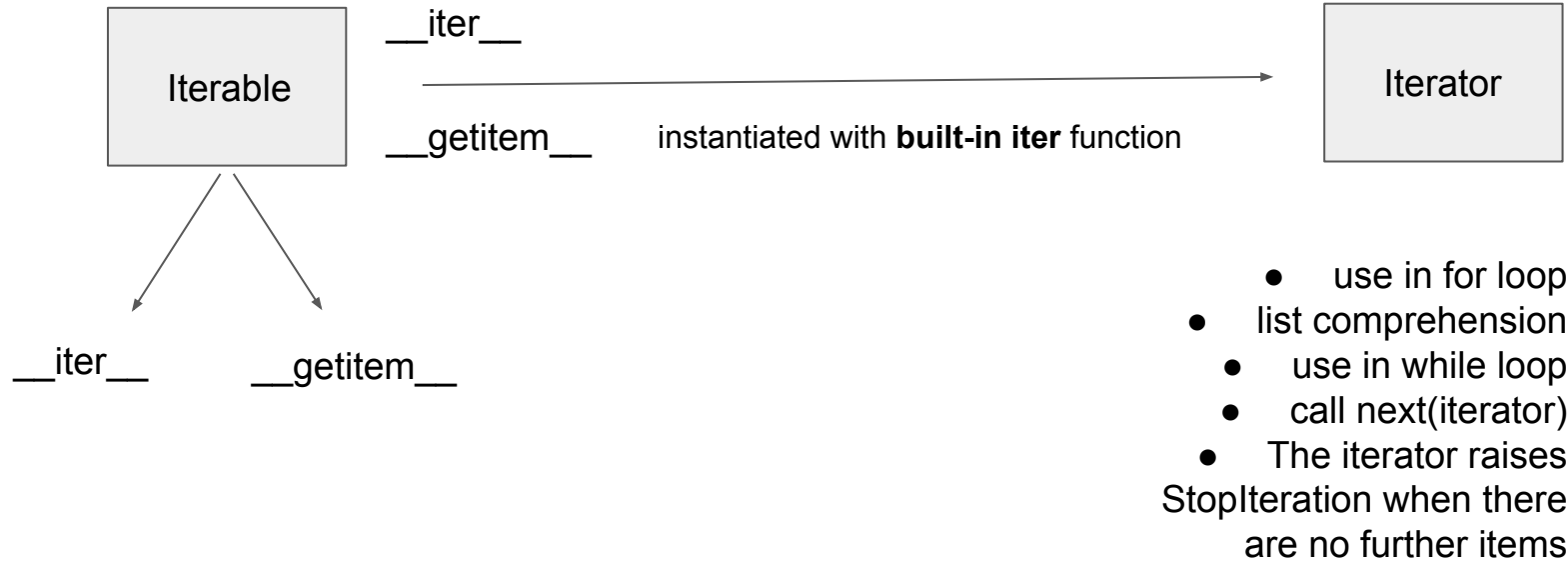
What is c?

Broadcasting can only be performed when the shape of each dimension in the arrays are equal or one has the dimension size of 1

Note: for broadcasting purpose, 1-d ndarray can be thought of as a row vector

Python Generators

- Iterable
- Iterators
- Generators



```
import re
import reprlib
```

```
RE_WORD = re.compile(r'\w+')
```

```
class Sentence:
```

```
    def __init__(self, text):
        self.text = text
        self.words = RE_WORD.findall(text)
```

```
    def __repr__(self):
        return f'Sentence({reprlib.repr(self.text)})'
```

```
    def __getitem__(self, index):
        return self.words[index]
```

```
    def __len__(self):
        return len(self.words)
```

```
s = Sentence("Hello world")
it = iter(s)
for i in it:
    print(i)
```

```
import re
import reprlib
```

```
RE_WORD = re.compile(r'\w+')
```

```
class Sentence:
```

```
    def __init__(self, text):
        self.text = text
        self.words = RE_WORD.findall(text)

    def __repr__(self):
        return f'Sentence({reprlib.repr(self.text)})'

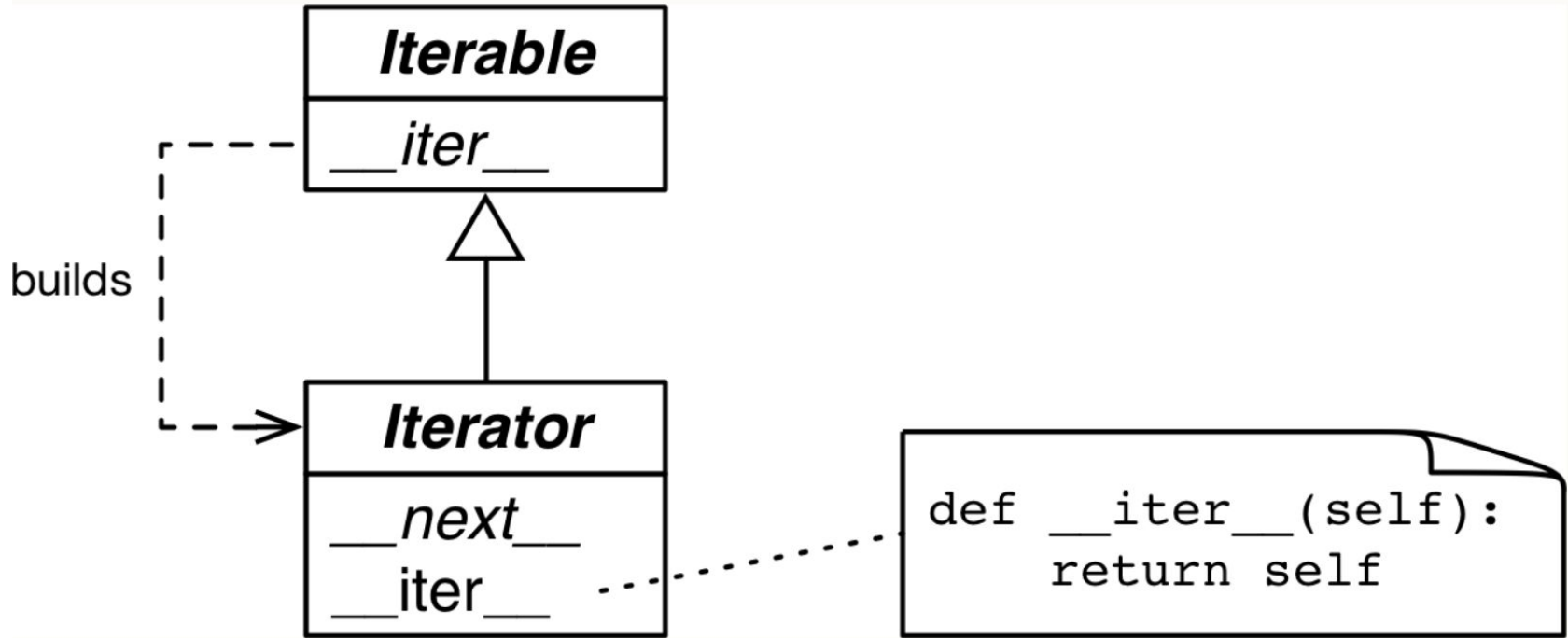
    def __iter__(self):
        return Sentenceliterator(self.words)
```

```
class Sentenceliterator:
```

```
    def __init__(self, words):
        self.words = words
        self.index = 0

    def __next__(self):
        try:
            word = self.words[self.index]
        except IndexError:
            raise StopIteration()
        self.index += 1
        return word

    def __iter__(self):
        return self
```



Therefore, iterators are also iterable, but iterables are not iterators.

```
import re
import reprlib

RE_WORD = re.compile(r'\w+')

class Sentence:

    def __init__(self, text):
        self.text = text
        self.words = RE_WORD.findall(text)

    def __repr__(self):
        return 'Sentence(%s)' % reprlib.repr(self.text)

    def __iter__(self):
        for word in self.words:
            yield word
```

Definition:

*a short-cut way to defining an iterator. All you need to do is define a function with at least 1 call to **yield** and now when you call that function it will return "**something**" which will act like an iterator (you can call **next** method and use it in a **for** loop)*


```
>>> def gen_123():
...     yield 1
...     yield 2
...     yield 3
...
>>> gen_123 # doctest: +ELLIPSIS
<function gen_123 at 0x...>
>>> gen_123() # doctest: +ELLIPSIS
<generator object gen_123 at 0x...>
>>> for i in gen_123():
...     print(i)
1
2
3
>>> g = gen_123()
>>> next(g)
1
>>> next(g)
2
>>> next(g)
3
>>> next(g)
Traceback (most recent call last):
...
StopIteration
```

```
>>> def gen_AB():
...     print('start')
...     yield 'A'
...     print('continue')
...     yield 'B'
...     print('end.')
...
>>> for c in gen_AB():
...     print('-->', c)
...
start
--> A
continue
--> B
end.
>>>
```



- Twitter API needs developer verification, so start early
- Announcement
 - Late submissions are not graded
 - One exception can be provided upon 48 hours prior notice

What is Pandas?

Wrapper around ndarray with more information such as column names, index etc

What is the difference between ndarray and list?

- ndarray - homogeneous, list - non-homogenous

What is the method to get the underlying ndarray of a pandas dataframe or series?

- `df.values()`

Data Structure

Series -> 1D

Dataframe -> 2D

Series

| | |
|--|--|
| <code>s = pd.Series(ndarray or dict, index=[], name=[])</code> | |
| | |
| <code>s[['idx1', 'idx2']]</code> | |
| <code>s[s > 0.5]</code> | |

Properties

| | |
|---------------------------|--|
| <code>s.index.name</code> | |
| <code>s.name</code> | |
| | |

Dataframe

| | |
|---|-----------------------------------|
| <code>df = pd.DataFrame(ndarray or dict, index=[], columns=[])</code> | |
| <code>df.head()</code> | |
| <code>df.tail()</code> | |
| <code>df.sample(10)</code> | |
| <code>df.loc(first arg: row index, sec arg: col names)</code> | |
| <code>df.loc[df['col1'] > 0.5]</code> | |
| <code>df.iloc(0, 1)</code> | |
| <code>apply(axis=[0 or 1])</code> | computation along rows or columns |
| <code>applymap()</code> | element-wise |

I/O

| | |
|-----------------------|--|
| read_csv / to_csv | |
| read_html / to_html | |
| read_excel / to_excel | |
| read_hdf / to_hdf | |
| many more... | |

General

| | |
|--------------------------|--|
| <code>df.describe</code> | |
| <code>df.info()</code> | |
| <code>df.isnull()</code> | |
| <code>df.dropna()</code> | |
| <code>df.fillna()</code> | |

Documentation is your friend!

<https://pandas.pydata.org/docs/reference/io.html>

Class exercise:

Write a python program to do matrix multiplication of two multi-dimensional without using np.dot and check your answer with it in the end:

```
a = [[1, 2, 3],  
      [4, 5, 6]]
```

```
b = [[7, 8],  
      [9, 10],  
      [11, 12]]
```

```
res = [[-1, -1],  
        [-1, -1]]
```

