

Key Term

- **numpy.ndarray** - n-dimensional Numpy array, the core object that underpins Pandas DataFrames
- **pyspark.sql.DataFrame** - Distributed dataframe object in PySpark for working with large datasets
- **dask.dataframe.DataFrame** - Dask dataframe that partitions Pandas dataframe across multiple cores
- **np.random.randn()** - Generate an ndarray with random values from a normal distribution
- **df.persist()** - Trigger caching of a PySpark dataframe into memory or disk storage

Pandas Code Examples

```
1 import pandas as pd
2 import numpy as np
3
4 # numpy.ndarray
5 data = np.random.randn(5, 4) # Generate 5x4 ndarray
6 df = pd.DataFrame(data) # Underlying structure of Pandas DataFrame
7 print(df)
8
9 # np.random.randn()
10 df = pd.DataFrame(np.random.randn(5, 4),
11 | | | | | columns=['A', 'B', 'C', 'D'])
12 print(df)
```

	0	1	2	3
0	-0.030742	0.477710	0.029938	-1.005135
1	-0.137934	0.299512	-1.404028	-0.153353
2	-0.755672	-0.451716	-0.182981	2.478212
3	-0.089476	0.472844	0.329931	-0.630641
4	-1.144661	-0.177849	-1.287844	-1.102989

	A	B	C	D
0	0.465819	-1.195653	0.618216	0.483000
1	-1.423888	0.671257	-1.298546	-1.954487
2	0.180366	-2.034974	-0.939537	-0.191529
3	-0.434719	0.248863	-2.398414	-0.791016
4	-0.191686	-0.278751	0.525708	-2.467836

Creating NumPy Arrays in Python

NumPy

- Designed for Scientific Computing
- Ndarray object = stand for N-dimensional array object → refer to as an array
- Highly optimized computational tools
- Ecosystem of Scientific and Mathematical Libraries

↳ Machine Learning: There are many the machine learning libraries that expect an array object as input

↳ if they don't, they will turn the input into an array object behind the scenes and return an array object as a result

ndarray

- Size set at creation
 - Single data type
 - Not limited to two dimensions
 - Can be reshaped returning a new array object
- Its size is set ↴ cannot resize it without creating a new object

NumPy Arrays

```
[65]: import numpy as np
```

Creating Arrays

```
[66]: data = [1, 2, 3, 4] ↴ first_array = np.array(data)  
first_array ↴ be an array, with that list  
[66]: array([1, 2, 3, 4])  
  
[67]: data = [[1, 2, 3],  
           [4, 5, 6],  
           [7, 8, 9]] ↴ my_array = np.array(data)  
my_array ↴ be an array with the list of lists  
[67]: array([[1, 2, 3],  
           [4, 5, 6],  
           [7, 8, 9]])
```

Creating

```
[69]: np.ones(12) ↴ make an array 1 twelve times  
[69]: array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.])  
  
[70]: np.zeros(3) ↴ also zeros  
[70]: array([0., 0., 0.])  
  
[71]: np.arange(10) ↴ create an array with list of range  
[71]: array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])  
      Starting number ↴ Ending number  
[72]: np.arange(3, 12, 4) ↴ step  
[72]: array([ 3,  7, 11])
```

Dimensions and Shape

```
[73]: oned = np.arange(21)  
oned  
[73]: array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13, 14, 15, 16,  
           17, 18, 19, 20])  
  
[74]: oned.shape  
[74]: (21,) ↴ cause it 1D  
  
[75]: oned.ndim  
[75]: 1 ↴ check its number of dimensions  
  
[76]: oned.size  
[76]: 21 ↴ list of lists  
[77]: list_of_lists = [[1,2,3], [4,5,6], [7,8,9]]  
twod = np.array(list_of_lists)  
twod  
[77]: array([[1, 2, 3],  
           [4, 5, 6],  
           [7, 8, 9]])  
  
[78]: twod.ndim  
[78]: 2  
  
[79]: twod.size  
[79]: 6  
  
[80]: twod.shape  
[80]: (3, 3)  
  
[81]: oned = np.arange(12)  
oned  
[81]: array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
```

↳ 3x6 = 18 correct
↳ column

```
[82]: twod = oned.reshape(3,6) ↴ Use reshape method to return array with different shape  
[82]: twod  
[82]: array([[ 0,  1,  2,  3],  
           [ 4,  5,  6,  7],  
           [ 8,  9, 10, 11]])  
  
[83]: twod.shape  
[83]: (3, 6)  
  
[84]: twod.ndim  
[84]: 2 ↴ from 1D incorrect to 2D  
[85]: twod.reshape(2,3) ↴ to use reshape, those arguments must be numbers that will add up to the right number  
ValueError: Traceback (most recent call last)  
/var/folders/29/5dyw2t0n71v13sp198c30g80000n/T/ipykernel_30144/198994001.py in <module>  
      1 twod.reshape(2,3) ↴ like how IF to turn shape 12 to shape 6  
ValueError: cannot reshape array of size 12 into shape (2,3)  
      2 ↴ 2x3=6  
      3 ↴ 2x2x3=12 correct  
  
[86]: twod.reshape(2,2,3)  
  
[86]: array([[[ 0,  1,  2],  
           [ 3,  4,  5]],  
           [[ 6,  7,  8],  
           [ 9, 10, 11]]])
```

Cont. Creating NumPy Arrays in Python

Setting Data Type

```
[87]: darray = np.arange(100)
darray
```

Setting Data Type → one of the things that separates NumPy from Pandas, is that they can only contain a single data type

```
[87]: 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, 32, 33,
34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50,
51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67,
68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84,
85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99])
```

```
[88]: darray.dtype
```

check its type

```
[88]: dtype('int64')
```

I can also check its size

```
[89]: darray.nbytes
```

```
[89]: 800
```

```
[90]: darray = np.arange(100, dtype=np.int8) # we can add new argument to specify a particular data type
darray
```

```
[90]: 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, 32, 33,
34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50,
51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67,
68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84,
85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99],
dtype='int8')
```

```
[91]: darray.dtype
```

```
[91]: dtype('int8')
```

```
[92]: darray.nbytes
```

```
[92]: 100 decrease memory
```

↳ we can fit them to small memory

- Once you set its datatype, NumPy will not let you add items of a different data type

```
[93]: darray[12] = 'a'
```

```
ValueError Traceback (most recent call last)
/var/folders/29/5dyw2t0n71v13sp1980c30g80000gr/T/ipykernel_30144/1692319446.py in <module>
      1 darray[12] = 'a'

ValueError: invalid literal for int() with base 10: 'a'
```

```
[94]: darray[12] = 0.4 ↳ its data type is the int
      ↳ but this value is decimal
```

```
[95]: darray[12]
```

so it stripped off

```
[95]: 0
```

↳ so be aware of data type, when using NumPy array

Broadcasting

Broadcasting

```
[96]: A1 = np.array([[1,2,3],
[4,5,6],
[7,8,9]])
```

```
[97]: A1.shape
```

```
[97]: (3, 3)
```

```
[ ]:
```

```
[98]: A2 = np.array([[1,1,1],
[1,1,1],
[1,1,1]])
```

```
[99]: A1 + A2 ↳ values within the arrays
are added together
```

```
[99]: array([[ 2,  3,  4],
[ 5,  6,  7],
[ 8,  9, 10]])
```

```
[100]: A1 + 3 ↳ we can add a single value
to the array
```

```
[100]: array([[ 4,  5,  6],
[ 7,  8,  9],
[10, 11, 12]])
```

```
[101]: A3 = np.array([1,1,1])
A3
```

```
[102]: A3.shape
```

```
[102]: (3,)
```

expanded out, it chose the maximum dimensions of the input arrays

```
[103]: A1 + A3 ↳ its still added to
      every single number
```

```
[103]: array([[ 2,  3,  4],
[ 5,  6,  7],
[ 8,  9, 10]]) ↳ and becomes
      a 3 by 3 array
```

This is called broadcasting

```
[104]: A4 = np.arange(10).reshape(2,5)
A4
```

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

```
[105]: A5 = np.arange(14).reshape(2,7,1)
A5
```

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

```
[106]: A6 = A4 + A5
```

```
[107]: A6.shape
```

```
[107]: (2, 7, 5)
```

Notice it selected by maximum size from each dimension

Cont. Creating NumPy Arrays in Python

Matrix Operations

```
[113]: M1 = np.arange(9).reshape(3,3)
M1
[113]: array([[0, 1, 2],
              [3, 4, 5],
              [6, 7, 8]])

[116]: M2 = np.arange(2, 11).reshape(3,3)
M2
[116]: array([[ 2,  3,  4],
              [ 5,  6,  7],
              [ 8,  9, 10]])

[117]: M1.transpose()
[117]: array([[0, 3, 6],
              [1, 4, 7],
              [2, 5, 8]])

[118]: M1.diagonal()
[118]: array([0, 4, 8])
          ↓ matrix product
[119]: M1 @ M2
[119]: array([[ 21,  24,  27],
              [ 66,  78,  90],
              [111, 132, 153]])
```

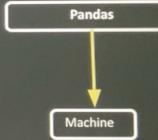
When to use them

ndarray	DataFrame
<ul style="list-style-type: none">• Multi dimensional data• Single <u>data type</u>• <u>Complex numerical calculations</u>	<ul style="list-style-type: none">• <u>Two dimensional data</u>• <u>Multiple data types</u>• <u>Data analysis</u>• <u>Data visualization</u>

Spark and PySpark DataFrames in Python

Pandas

- Designed to run on single machine
- Performance bound by machines memory
- Chunking → load a subset of data and work on that and then load another chunk of data
- Upper limit Gigabytes (1? 5? 100?)
→ generally



Big Data → Petabytes, Exabytes

- Hadoop
- Spark
- Distributed computing

These nodes could be different computers, different virtual machines, or different sub-processes

Node 1

Node 1

Node 1

Node 1

Node 1

uses multiple nodes

allow to scale horizontally

add machines to increase performance not necessarily have to add much larger and more expensive machines

Hadoop: does this by writing intermediate files during calculations

↳ doesn't require a great amount of RAM

Spark:

- newer solution
- more performant, does work in memory.
- by using memory all nodes

Spark

- Distributed DataFrames on JVM
↳ Java Virtual machine
- Written in Scala
- PySpark Library → include, use spark DataFrame by using Python context
- Data sources include Hadoop HDFS, S3, and Streaming
- Lazy Evaluation

- its core uses DataFrames ← Distributed data frame

Eager (Pandas)



↳ every operation in Pandas code calculates its result before the next operation is run
↳ useful when interactively, intuitive to debug

Lazy (Spark)

↳ taking advantage of PySpark

↳ Operations are stacked and optimized transformation is calculated by Spark ← run in background
↳ use its nature distributed

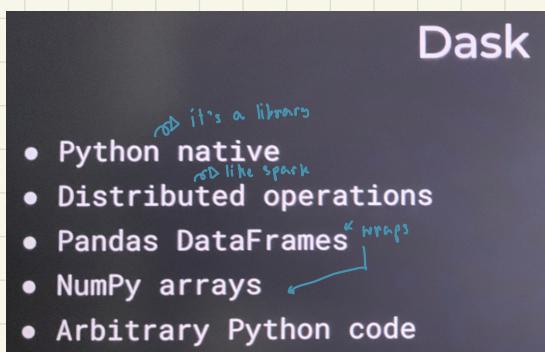
Considerations

between Pandas or Spark

→ more than Gb, use spark
↳ <= Gb like csv file, use pandas

- Built for big data transformations
↳ are you plan for the way you write your transformation? Eager or Lazy?
- Optimized integration plan strategy
↳ Spark is challenging to do a debugging sometime
- Debugging challenges

Creating Dash DataFrames in Python



```
[11]: import dask.dataframe as dd
[12]: df = pd.DataFrame(data)
       ↓ method → data → num. that we wanted to
       divide the data
[13]: ddf = dd.from_pandas(df, npartitions=3)
[13]: Dask DataFrame Structure:
      a   b
npartitions=3
      0   int64  int64
      33333334 ... ...
      66666668 ... ...
      99999999 ... ...
      ... ...
      no actual values
Dask Name: from_pandas, 3 tasks
```

Task

```
[1]: import pandas as pd
import random

[4]: leng = 1000000000 ← large now
data = {'a': (random.randint(0, 100) for _ in range(leng)),
        'b': (random.randint(2, 200) for _ in range(leng))}

[5]: df = pd.DataFrame(data)
df.head()

[5]:
```

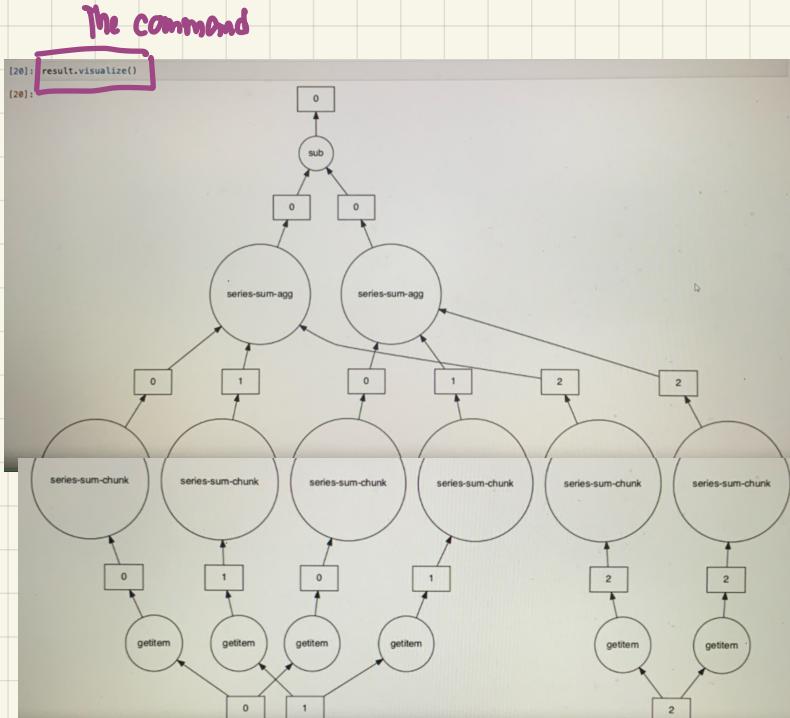
a	b
0	94 130
1	0 167
2	41 40
3	97 146
4	46 6

```
[10]: df.std()
before complete this,
It took time for a while
[10]: a    29.153448
      b    57.449703
      dtype: float64
```

```
[14]: ddf.std() → It didn't give
[14]: Dask Series Structure:
      npartitions=1
      a    float64
      b    ...
      dtype: float64
      Dask Name: dataframe-std, 15 tasks
[15]: ddf.std().compute() ← To force Dask to
      execute
[15]: a    29.153448
      b    57.449703
      dtype: float64
[16]: result = ddf.a.sum() - ddf.b.sum()
      result
[16]: dd.Scalar<sub-0ad..., dtype=int64>
[17]: result.compute()
[17]: -5099626371
```

```
[18]: result.task
[18]: HighLevelGraph
HighLevelGraph with 8 layers.
● Layer1: from_pandas
○ Layer2: getitem
○ Layer3: series-sum-chunk
● Layer4: series-sum-agg
○ Layer5: getitem
○ Layer6: series-sum-chunk
● Layer7: series-sum-agg
● Layer8: sub
```

It creates a graph behind the scenes



Choosing a framework

Pandas or NumPy

small-medium Size of data

Pandas

visualization

2 dim
various data type

PySpark

Vs

PySpark or Dask

large

Vs

NumPy

more dim
single data type

Dask

large ecosystem Nature of enterprise Python native

all pipeline with Python ↗