

Pandas

Data Manipulation in Python

Pandas

- ▶ Built on NumPy
- ▶ Adds data structures and data manipulation tools
- ▶ Enables easier data cleaning and analysis

```
1 import pandas as pd
2 pd.set_option("display.width", 120)
```

That last line allows you to display DataFrames with many columns without wrapping.

Pandas Fundamentals

Three fundamental Pandas data structures:

- ▶ **Series** - a one-dimensional array of values indexed by a `pd.Index`
- ▶ **Index** - an array-like object used to access elements of a **Series** or **DataFrame**
- ▶ **DataFrame** - a two-dimensional array with flexible row indices and column names

Series from List

```
1 In [4]: data = pd.Series(['a', 'b', 'c', 'd'])
2
3 In [5]: data
4 Out[5]:
5 0    a
6 1    b
7 2    c
8 3    d
9 dtype: object
```

The 0..3 in the left column are the `pd.Index` for `data`:

```
1 In [7]: data.index
2 Out[7]: RangeIndex(start=0, stop=4, step=1)
```

The elements from the Python list we passed to the `pd.Series` constructor make up the values:

```
1 In [8]: data.values
2 Out[8]: array(['a', 'b', 'c', 'd'], dtype=object)
```

Notice that the values are stored in a Numpy array.

Series from Sequence

You can construct a list from any definite sequence:

```
1 In [24]: pd.Series(np.loadtxt('exam1grades.txt'))
2 Out[24]:
3 0      72.0
4 1      72.0
5 2      50.0
6 ...
7 134    87.0
8 dtype: float64
```

or

```
1 In [25]: pd.Series(open('exam1grades.txt').readlines())
2 Out[25]:
3 0      72\n
4 1      72\n
5 2      50\n
6 ...
7 134    87\n
8 dtype: object
```

... but not an indefinite sequence:

```
1 In [26]: pd.Series(open('exam1grades.txt'))
2 ...
3 TypeError: object of type '_io.TextIOWrapper' has no len()
```

Series from Dictionary

```
1 salary = {"Data Scientist": 110000,  
2           "DevOps Engineer": 110000,  
3           "Data Engineer": 106000,  
4           "Analytics Manager": 112000,  
5           "Database Administrator": 93000,  
6           "Software Architect": 125000,  
7           "Software Engineer": 101000,  
8           "Supply Chain Manager": 100000}
```

Create a `pd.Series` from a `dict`:

```
1 In [14]: salary_data = pd.Series(salary)  
2  
3 In [15]: salary_data  
4 Out[15]:  
5 Analytics Manager      112000  
6 Data Engineer          106000  
7 Data Scientist         110000  
8 Database Administrator  93000  
9 DevOps Engineer        110000  
10 Software Architect     125000  
11 Software Engineer     101000  
12 Supply Chain Manager   100000  
13 dtype: int64
```

The index is a sorted sequence of the keys of the dictionary passed

Series with Custom Index

General form of Series constructor is `pd.Series(data, index=index)`

- ▶ Default is integer sequence for sequence data and sorted keys of dictionaries
- ▶ Can provide a custom index:

```
1 In [29]: pd.Series([1,2,3], index=['a', 'b', 'c'])
2 Out[29]:
3 a      1
4 b      2
5 c      3
6 dtype: int64
```

The index object itself is an immutable array with set operations.

```
1 In [30]: i1 = pd.Index([1,2,3,4])
2
3 In [31]: i2 = pd.Index([3,4,5,6])
4
5 In [32]: i1[1:3]
6 Out[32]: Int64Index([2, 3], dtype='int64')
7
8 In [33]: i1 & i2 # intersection
9 Out[33]: Int64Index([3, 4], dtype='int64')
10
11 In [34]: i1 | i2 # union
```

Series Indexing and Slicing

Indexing feels like dictionary access due to flexible index objects:

```
1 In [37]: data = pd.Series(['a', 'b', 'c', 'd'])
2
3 In [38]: data[0]
4 Out[38]: 'a'
5
6 In [39]: salary_data['Software Engineer']
7 Out[39]: 101000
```

But you can also slice using these flexible indices:

```
1 In [40]: salary_data['Data Scientist':'Software Engineer']
2 Out[40]:
3 Data Scientist          110000
4 Database Administrator  93000
5 DevOps Engineer        110000
6 Software Architect     125000
7 Software Engineer      101000
8 dtype: int64
```


Basic DataFrame Structure

A DataFrame is a series Serieses with the same keys. For example, consider the following dictionary of dictionaries meant to leverage your experience with spreadsheets (in [spreadsheet.py](#)):

```
1 In [5]: import spreadsheet; spreadsheet.cells
2
3 Out[5]:
4 {'A': {1: 'A1', 2: 'A2', 3: 'A3'},
5  'B': {1: 'B1', 2: 'B2', 3: 'B3'},
6  'C': {1: 'C1', 2: 'C2', 3: 'C3'},
7  'D': {1: 'D1', 2: 'D2', 3: 'D3'}}
```

All of these dictionaries have the same keys, so we can pass this dictionary of dictionaries to the DataFrame constructor:

```
1 In [7]: ss = pd.DataFrame(spreadsheet.cells); ss
2
3 Out[7]:
4      A  B  C  D
5 1  A1 B1 C1 D1
6 2  A2 B2 C2 D2
7 3  A3 B3 C3 D3
```

Basic DataFrame Structure

```
1 In [5]: import spreadsheet; spreadsheet.cells
2
3 Out[5]:
4 {'A': {1: 'A1', 2: 'A2', 3: 'A3'},
5  'B': {1: 'B1', 2: 'B2', 3: 'B3'},
6  'C': {1: 'C1', 2: 'C2', 3: 'C3'},
7  'D': {1: 'D1', 2: 'D2', 3: 'D3'}}
```

All of these dictionaries have the same keys, so we can pass this dictionary of dictionaries to the DataFrame constructor:

```
1 In [7]: ss = pd.DataFrame(spreadsheet.cells); ss
2
3 Out[7]:
4      A  B  C  D
5 1  A1 B1 C1 D1
6 2  A2 B2 C2 D2
7 3  A3 B3 C3 D3
```

- ▶ Each column is a Series whose keys (index) are the values printed to the left (1, 2 and 3).
- ▶ Each row is a Series whose keys (index) are the column headers.

Try evaluating `ss.columns` and `ss.index`.

DataFrame Example

Download [hotjobs.py](#) and do a `%load hotjobs.py` (to evaluate the code in the top-level namespace instead of importing it).

```
1 In [42]: jobs = pd.DataFrame({'salary': salary, 'openings': openings})
2
3 In [43]: jobs
4 Out[43]:
```

	openings	salary
6 Analytics Manager	1958	112000
7 Data Engineer	2599	106000
8 Data Scientist	4184	110000
9 Database Administrator	2877	93000
10 DevOps Engineer	2725	110000
11 Software Architect	2232	125000
12 Software Engineer	17085	101000
13 Supply Chain Manager	1270	100000
14 UX Designer	1691	92500

```
1 In [46]: jobs.index
2 Out[46]:
3 Index(['Analytics Manager', 'Data Engineer', 'Data Scientist',
4       'Database Administrator', 'DevOps Engineer', 'Software Architect',
5       'Software Engineer', 'Supply Chain Manager', 'UX Designer'],
6       dtype='object')
7
8 In [47]: jobs.columns
```

Simple DataFrame Indexing

Simplest indexing of DataFrame is by column name.

```
1 In [48]: jobs['salary']
2 Out[48]:
3 Analytics Manager      112000
4 Data Engineer         106000
5 Data Scientist        110000
6 Database Administrator 93000
7 DevOps Engineer       110000
8 Software Architect    125000
9 Software Engineer     101000
10 Supply Chain Manager  100000
11 UX Designer           92500
12 Name: salary, dtype: int64
```

Each column is a Series:

```
1 In [49]: type(jobs['salary'])
2 Out[49]: pandas.core.series.Series
```

General Row Indexing

The `loc` indexer indexes by row name:

```
1 In [13]: jobs.loc['Software Engineer']
2 Out[13]:
3 openings      17085
4 salary        101000
5 Name: Software Engineer, dtype: int64
6
7 In [14]: jobs.loc['Data Engineer':'Database Administrator']
8 Out[14]:
9                openings salary
10 Data Engineer      2599 106000
11 Data Scientist     4184 110000
12 Database Administrator 2877  93000
```

Note that slice ending is inclusive when indexing by name.

The `iloc` indexer indexes rows by position:

```
1 In [15]: jobs.iloc[1:3]
2 Out[15]:
3                openings salary
4 Data Engineer      2599 106000
5 Data Scientist     4184 110000
```

Note that slice ending is exclusive when indexing by integer position.

Special Case Row Indexing

```
1 In [16]: jobs[:2]
2 Out[16]:
3           openings salary
4 Analytics Manager    1958 112000
5 Data Engineer       2599 106000
6
7 In [17]: jobs[jobs['salary'] > 100000]
8 Out[17]:
9           openings salary
10 Analytics Manager    1958 112000
11 Data Engineer       2599 106000
12 Data Scientist      4184 110000
13 DevOps Engineer     2725 110000
14 Software Architect  2232 125000
15 Software Engineer   17085 101000
```

Try `jobs['salary'] > 100000` by itself. What's happening in `In[17]` above?

loc and iloc Indexing

The previous examples are shortcuts for `loc` and `iloc` indexing:

```
1 In [20]: jobs.iloc[:2]
2 Out[20]:
3           openings salary
4 Analytics Manager    1958 112000
5 Data Engineer       2599 106000
6
7 In [21]: jobs.loc[jobs['salary'] > 100000]
8 Out[21]:
9           openings salary
10 Analytics Manager    1958 112000
11 Data Engineer       2599 106000
12 Data Scientist      4184 110000
13 DevOps Engineer     2725 110000
14 Software Architect  2232 125000
15 Software Engineer   17085 101000
```

Aggregate Functions

The values in a series is a `numpy.ndarray`, so you can use NumPy functions, broadcasting, etc.

▶ Average salary for all these jobs:

```
1 In [14]: np.average(jobs['salary'])
2 Out[14]: 107125.0
```

▶ Total number of openings:

```
1 In [15]: np.sum(jobs['openings'])
2 Out[15]: 34930
```

And so on.

Adding Columns by Broadcasting

Add column by broadcasting a constant value:

```
1 In [16]: jobs['DM Prepares'] = True
2
3 In [17]: jobs
4 Out[17]:
```

	openings	salary	DM Prepares
6 Analytics Manager	1958	112000	True
7 Data Engineer	2599	106000	True
8 Data Scientist	4184	110000	True
9 Database Administrator	2877	93000	True
10 DevOps Engineer	2725	110000	True
11 Software Architect	2232	125000	True
12 Software Engineer	17085	101000	True
13 Supply Chain Manager	1270	100000	True

Adding Column by Applying Ufuncs

```
1 In [25]: jobs['Percent Openings'] = jobs['openings'] /
      np.sum(jobs['openings'])
2
3 In [26]: jobs
4 Out[26]:
```

	openings	salary	DM Prepares	Percent Openings
6 Analytics Manager	1958	112000	True	0.056055
7 Data Engineer	2599	106000	True	0.074406
8 Data Scientist	4184	110000	True	0.119782
9 Database Administrator	2877	93000	True	0.082365
10 DevOps Engineer	2725	110000	True	0.078013
11 Software Architect	2232	125000	True	0.063899
12 Software Engineer	17085	101000	True	0.489121
13 Supply Chain Manager	1270	100000	True	0.036358

CSV Files

Pandas has a very powerful CSV reader. Do this in iPython (or `help(pd.read_csv)` in the Python REPL):

```
1 pd.read_csv?
```

Now let's read the `super-grades.csv` file and re-do [Calc Grades](#) exercise using Pandas.

- ▶ Note that there is a missing Exam 2 grade for Farva.

Read a CSV File into a DataFrame

`super-grades.csv` contains:

```
1 Student,Exam 1,Exam 2,Exam 3
2 Thorny,100,90,80
3 Mac,88,99,111
4 Farva,45,,67
5 Rabbit,59,61,67
6 Ursula,73,79,83
7 Foster,89,97,101
```

First line is header, which Pandas will infer, and we want to use the first column for index values by passing a value to the `index_col` parameter:

```
1 In [3]: sgs = pd.read_csv('super-grades.csv', index_col=0)
2
3 In [4]: sgs
4 Out[4]:
5      Exam 1 Exam 2 Exam 3
6 Student
7 Thorny    100    90    80
8 Mac      88    99   111
9 Farva    45    56    67
10 Rabbit   59    61    67
11 Ursula   73    79    83
12 Foster   89    97   101
```

Applying an Arbitrary Function to Values in Rows

Now let's add a column with the average grades for each student by applying a `course_avg` function.

```
1 def course_avg(row):  
2     # Drop lowest grade  
3     return np.mean(row.values)
```

If we apply this to the DataFrame we get a Series with averages. Notice that we're "collapsing" columns (`axis=1`), that is, calculating values from a row like we did in NumPy:

```
1 In [6]: sgs.apply(course_avg, axis=1)  
2 Out[6]:  
3 Student  
4 Thorny    95.0  
5 Mac      105.0  
6 Farva    NaN  
7 Rabbit   64.0  
8 Ursula   81.0  
9 Foster   99.0  
10 dtype: float64
```

Before we add this series as a new column to our DataFrame we need to deal with the missing value for Farva.

Dealing with Missing Values

Many approaches. A simple one: replace NaN with a particular value:

```
1 In [8]: sgs.fillna(0)
2 Out[8]:
3         Exam 1 Exam 2 Exam 3
4 Student
5 Thorny      100   90.0    80
6 Mac         88   99.0   111
7 Farva       45    0.0    67
8 Rabbit      59   61.0    67
9 Ursula      73   79.0    83
10 Foster     89   97.0   101
```

A better approach: use a nan-ignoring aggregate function:

```
1 def course_avg_excuse_missing(row):
2     return np.nanmean(np.sort(row.values)[1:])
```

Adding Columns Calculated by Arbitrary Functions

Applying a function to rows creates a Series that we add to the DataFrame:

```
1 In [14]: sgs["avg"] = sgs.apply(course_avg_excuse_missing, axis=1); sgs
2 Out[14]:
3           Exam 1 Exam 2 Exam 3          avg
4 Student
5 Thorny         100   90.0    80   91.250000
6 Mac           88   99.0   111  100.750000
7 Farva         45    NaN    67   59.666667
8 Rabbit        59   61.0    67   62.750000
9 Ursula        73   79.0    83   79.000000
10 Foster       89   97.0   101   96.500000
```

Appending DataFrames

Now let's add a new row containing the averages for each exam.

- ▶ We can get the item averages by applying `np.mean` to the columns (`axis=0` – “collapsing” rows):

```
1 In [35]: sgs.apply(np.mean, axis=0)
2 Out[35]:
3 Exam 1    75.666667
4 Exam 2    80.333333
5 Exam 3    84.833333
6 avg      80.277778
7 dtype: float64
```

- ▶ We can turn this Series into a DataFrame with the label we want:

```
1 In [38]: pd.DataFrame({"ItemAverage": sgs.apply(np.mean, axis=0)})
2 Out[38]:
3           ItemAverage
4 Exam 1    75.666667
5 Exam 2    80.333333
6 Exam 3    84.833333
7 avg      80.277778
```


DataFrame Transpose

We need to give ItemAverages the same shape as our grades DataFrame:

```
1 In [41]: item_avgs = pd.DataFrame({"Item Avg": sgs.apply(np.mean,
2         axis=0)}).transpose()
3 In [43]: item_avgs
4 Out[43]:
5         Exam 1    Exam 2    Exam 3    avg
6 Item Avg 75.666667 80.333333 84.833333 80.277778
```

Then we can append the DataFrame because it has the same columns:

```
1 In [24]: sgs = sgs.append(item_avgs)
2
3 In [25]: sgs
4 Out[25]:
5         Exam 1    Exam 2    Exam 3    avg
6 Thorny 100.000000 90.000000 80.000000 90.000000
7 Mac   88.000000 99.000000 111.000000 99.333333
8 Farva 45.000000 56.000000 67.000000 56.000000
9 Rabbit 59.000000 61.000000 67.000000 62.333333
10 Ursula 73.000000 79.000000 83.000000 78.333333
11 Foster 89.000000 97.000000 101.000000 95.666667
12 Item Avg 75.666667 80.333333 84.833333 80.277778
```

Adding a Letter Grades Column

Adding a column with letter grades is easier than adding a column with a more complex calculation.

```
1 In [40]: sgs['Grade'] = \  
2     ...:     np.where(sgs['avg'] >= 90, 'A',  
3     ...:             np.where(sgs['avg'] >= 80, 'B',  
4     ...:                 np.where(sgs['avg'] >= 70, 'C',  
5     ...:                     np.where(sgs['avg'] >= 60, 'D',  
6     ...:                         'D'))))  
7     ...:  
8  
9 In [41]: sgs  
10 Out[41]:  
11          Exam 1    Exam 2    Exam 3    avg Grade  
12 Thorny    100.000000  90.000000  80.000000  90.000000    A  
13 Mac      88.000000  99.000000  111.000000 99.333333    A  
14 Farva    45.000000  56.000000  67.000000  56.000000    D  
15 Rabbit   59.000000  61.000000  67.000000  62.333333    D  
16 Ursula   73.000000  79.000000  83.000000  78.333333    C  
17 Foster   89.000000  97.000000  101.000000 95.666667    A  
18 Item Avg 75.666667  80.333333  84.833333  80.277778    B
```

Grouping and Aggregation

Grouping and aggregation can be conceptualized as a *split, apply, combine* pipeline.

► Split by Grade

1		Exam 1	Exam 2	Exam 3	avg Grade	
2	Thorny	100.000000	90.000000	80.000000	90.000000	A
3	Mac	88.000000	99.000000	111.000000	99.333333	A
4	Foster	89.000000	97.000000	101.000000	95.666667	A
1		Exam 1	Exam 2	Exam 3	avg Grade	
2	Item Avg	75.666667	80.333333	84.833333	80.277778	B
1		Exam 1	Exam 2	Exam 3	avg Grade	
2	Ursula	73.000000	79.000000	83.000000	78.333333	C
1		Exam 1	Exam 2	Exam 3	avg Grade	
2	Farva	45.000000	56.000000	67.000000	56.000000	D
3	Rabbit	59.000000	61.000000	67.000000	62.333333	D

- Apply some aggregation function to each group, such as sum, mean, count.
- Combine results of function applications to get final results for each group.

Letter Grades Counts

Here's how to find the counts of letter grades for our super troopers:

```
1 In [58]: sgs['Grade'].groupby(sgs['Grade']).count()
2 Out[58]:
3 Grade
4 A      3
5 B      1
6 C      1
7 D      2
8 Name: Grade, dtype: int64
```

Messy CSV Files

Remember the [Tides Exercise](#)? Pandas's `read_csv` can handle most of the data pre-processing:

```
1 pd.read_csv('wpb-tides-2017.txt', sep='\t', skiprows=14, header=None,
2             usecols=[0,1,2,3,5,7],
3             names=['Date', 'Day', 'Time', 'Pred(ft)', 'Pred(cm)',
4                   'High/Low'],
5             parse_dates=['Date', 'Time'])
```

Let's use the indexing and data selection techniques we've learned to re-do the [Tides Exercise](#) as a Jupyter Notebook. For convenience, `wpb-tides-2017.txt` is in the [code/analytics](#) directory, or you can [download it](#).

Reading SQL Databases

Let's create a realistically sized grades example dataset using fake student names. We'll get the names from the [Sakila Sample Database](#).

```
1 import numpy as np
2 import pandas as pd
3 import pymysql
4
5 sakila = pymysql.connect(host="localhost",
6                          user="root",
7                          password="",
8                          db="sakila",
9                          charset="utf8mb4",
10                         cursorclass=pymysql.cursors.DictCursor)
11
12 names = pd.read_sql("select first_name, last_name from actor", con =
13                    sakila)
14 names.head()
```

```
1      first_name  last_name
2  0  PENELOPE      GUINNESS
3  1      NICK      WAHLBERG
4  2      ED      CHASE
5  3  JENNIFER      DAVIS
6  4  JOHNNY      LOLLOBRIGIDA
```

Creating Datasets

Look at the [sakil-grades.ipynb](#) notebook for an example of extracting data from a database and creating a realistically sized fake data set similar to the grades file downloaded from Canvas.

There's also an [org-mode](#) version for Emacs users: [sakila-grades.org](#)