Data Processing and Analysis in Python Lecture 14 Tabular Data and pandas



DR. ADAM LEE

Pandas



- Panel Data System
- Built for the Python programming language
- An open source, BSD-licensed library (module)
- High-performance, easy-to-use data structures and data analysis tools
- Key components:
 - Series
 - DataFrame



Use Pandas

```
Import the whole module
>>> import pandas
  Import the whole module using alias
>>> import pandas as pd
# Import the whole module with all libraries
>>> from pandas import *
http://pandas.pydata.org/
>>> dir (pandas)
>>> dir (pandas. Series)
>>> dir (pandas. DataFrame)
>>> help(pandas)
>>> help (pandas. Series)
>>> help(pandas.DataFrame)
```



Lee 704 Pandas

Data Type: Series

- One-dimensional array like object containing data and labels (or index)
- Lots of ways to build a Series:

```
>>> srs=pd.Series(['a','b','c','d','e','f'])
>>> srs
0     a
1     b
2     c
3     d
4     e
5     f
dtype: object
```



Series from Dictionary

■ A dict-like container:

```
>>> dct={'a':100,'b':150,'c':200}
>>> srs=pd.Series(dct)
>>> srs
a    100
b    150
c    200
dtype: int64
>>> ses['b']
150
```



Series with Index

- A Series index can be specified:
 - Single values can be selected by index
 - Multiple values can be selected with multiple indexes

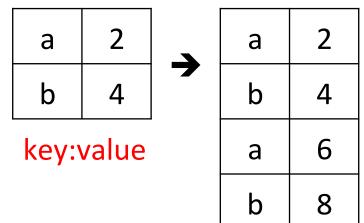


Series versus Dictionary

- Think of a Series as a fixed-length order dictionary
- However, unlike dictionary, index items do not have to be unique:
 dict

а	2		а	6
b	4	→	b	8

pandas.Series



Data: DataFrame

- Tabular data
- Rectangular data structure, a lot like an array
- Columns (Series) of different types
- Rows and columns act differently
- Can index by (column) labels as well as positions
- Handles missing data (NaN)
- Convenient plotting
- Fast operations with keys
- Lots of facilities for input/output



Series to DataFrame

Series

 pop

 0
 18.8

 1
 19.1

 2
 9.7

 3
 9.7

 4
 9.8

Series

	state
0	FL
1	FL
2	GA
3	GA
4	GA

Series

	year
0	2010
1	2011
2	2008
3	2010
4	2011

DataFrame

	pop	state	year
0	18.8	FL	2010
1	19.1	FL	2011
2	9.7	GA	2008
3	9.7	GA	2010
4	9.8	GA	2011

index columns (key) (values)

4	A		В			
1	рор	∇	state	∇	year	$\overline{\mathbf{v}}$
2	18	3.8	FL			2010
3	19	9.1	FL			2011
4	9	9.7	GA			2008
5	9	9.7	GA			2010
6	9	9.8	GA			2011

Create Data

Creation with dict of equal-length lists:

```
>>> data={"state":["FL","FL","GA","GA","GA"],
         "year" : [2010,2011,2008,2010,2011],
         "pop" : [18.8,19.1, 9.7, 9.7, 9.8]}
>>> frame=pd.DataFrame(data)
>>> frame
   pop state year
 18.8 FL 2010
 19.1 FL 2011
2 9.7 GA 2008
3 9.7 GA 2010
 9.8
      GA 2011
```



Lee 704 Pandas

Create Data

Creation with dict of dicts:

```
>>> pop data={"FL":{2010:18.8,
                    2011:19.1},
              "GA":{2008: 9.7,
                    2010: 9.7,
                    2011: 9.8}}
>>> pop=pd.DataFrame(pop data)
>>> pop
        TT.
            GΑ
2008 NaN 9.7
2010 18.8 9.7
2011 19.1 9.8
```



Lee 704 Pandas

10 ROBERT H. SMITH SCHOOL OF BUSINESS

Create List

The initial set of baby names and birth rates:

```
>>> names = ["Bob", "Jane", "Mary", "John", "Mel"]
>>> births = [968, 155, 77, 578, 973]
>>> BabyDataSet = list(zip(names, births))
>>> BabyDataSet
[('Bob', 968), ('Jane', 155), ('Mary', 77),
('John', 578), ('Mel', 973)]
```



Lee 704 Pandas

11 ROBERT H. SMITH

Create Data

Create a Pandas DetaFrame object df:

```
>>> df = pd.DataFrame(data=BabyDataSet,
                    columns=["names", "births"])
>>> df
        births
  names
    Bob
             968
0
             155
   Jane
             77
  Mary
3
   John
            578
             973
    Me l
```

 You can think of df holding the contents of BabyDataSet in a format similar to an SQL table or an Excel spreadshee

MARYLAND

Lee 704 Pandas

12 ROBERT H. SMI

Save Data into File

■ to_csv() Export DataFrame to a comma-separated values (CSV) file:

```
>>> df.to_csv("filename.csv")
>>> df.to_csv("filename.csv", index=False,
header=False)
```

The file will be saved in the same location as your python file unless specified otherwise



Lee 704 Pandas

13 ROBERT H. SMITH

Load Data from File

■ read csv() To import from the CSV file:



Lee 704 Pandas

14 ROBERT H. SMITH SCHOOL OF BUSINES.

Load Data from File

■ To give columns specific names, set parameter "names":

```
>>> filepath = "filename.csv"
>>> df = pd.read_csv(filepath, names=["names",
"births"])
>>> df
  names births
0  Bob  968
1  Jane  155
2  Mary  77
3  John  578
4  Mel  973
```



Lee 704 Pandas 15 ROBERT H. SMITH

Select Data

- df.head() returns the first 5 records
- df.head(n) returns the first n records
- df.tail() returns the last 5 records
- df.tail(n) returns the last n records



Lee 704 Pandas 16 ROBERT H. SMITH

Prepare Data

■ To find all unique records of "names" column:

```
>>> df["names"].unique()
array(['Bob', 'Jane', 'Mary', 'John', 'Mel'],
dtype=object)
```

■ To obtain all descriptive statistics:

```
>>> df["names"].describe()
count 5
unique 5
top Jane
freq 1
Name: Names, dtype: object
```



Lee 704 Pandas 17 ROBERT H. SMITH

Analyze Data

To find baby name with the highest birth rate:

```
>>> sorted = df.sort values(["births"],
ascending=False)
>>> sorted
 names births
   Mel
           973
4
 Bob 968
3 John 578
1 Jane
           155
            77
  Mary
>>> sorted.head(1)
 names births
4
   Mel
           973
```



Analyze Data

- df ["names"] returns the entire list of baby names, i.e. the entire "Names" column
- df["births"] returns the entire list of births, i.e. the entire
 "Births" column
- df["births"].max() returns the maximum value found in the "Births" column
- [df["births"] == df["births"].max()] equals to
 [Find all records in the "Births" column, where it equals to the
 maximum, i.e. 973.]
- df["names"][df["births"] ==
 df["births"].max()] equals to [Find all records in the
 "Names" column, where the "Births" column equals to the
 maximum, i.e. 973.]

Lee 704 Pandas 19 ROBERT H. SMITH

```
>>> data = range(0,10)
>>> df = pd.DataFrame(list(data))
>>> df.columns = ["num"]
>>> df["new"] = 5
>>> df["new"] = df["new"] + 1
```

	num	new
0	0	6
1	1	6
2	2	6
3	3	6
4	4	6
5	5	6
6	6	6
7	7	6
8	8	6
9	9	6

Lee 704 Pandas 20 ROBERT H. SMITH

```
>>> data = range(0,10)
>>> df = pd.DataFrame(list(data))
>>> df.columns = ["num"]
>>> df["three"] = 3
>>> df["col"] = df["num"]
```

	num	three	col
0	0	3	0
1	1	3	1
2	2	3	2
3	3	3	3
4	4	3	4
5	5	3	5
6	6	3	6
7	7	3	7
8	8	3	8
9	9	3	9

Lee 704 Pandas

21 ROBERT H. SMITH

	num	three	col
а	0	3	0
b	1	3	1
С	2	3	2
d	3	3	3
е	4	3	4
f	5	3	5
g	6	3	6
h	7	3	7
i	8	3	8
j	9	3	9

Lee 704 Pandas

22 ROBERT H. SMITH

```
>>> df.loc['a']
Rev 0
test 3
col 0
Name: a, dtype: int64
# loc[] both inclusive
>>> df.loc['a':'d']
# iloc[] inclusive:exclusive
>>> df.iloc[0:3]
```

	num	three	col
а	0	3	0
b	1	3	1
C	2	3	2
d	3	3	3

	num	three	col
а	0	3	0
b	1	3	1
С	2	3	2



```
# ...
>>> df[["num","three"]]

# df.ix[rows, columns]
>>> df.ix[0:3,"num"]
>>> df.ix[5:,"col"]
>>> df.ix[5:,"col"]
```

	num	three
а	0	3
b	1	3
С	2	3

а	0
b	1
С	2

f	5
g	6
h	7
.	8
j	9

	num	three
а	0	3
b	1	3
С	2	3
d	3	3
е	4	3
f	5	3
g	6	3
h	7	3
i	8	3
j	9	3

GroupBy Operations

	one	two
letter		
а	2	4
b	2	4
С	1	2

	letter	one	two
0	а	1	2
1	а	1	2
2	b	1	2
3	b	1	2
4	С	1	2

>>> letter.sum()

GroupBy Operations

```
>>> letter one =
df.groupby(["letter", "one"]).sum()
>>> letter one =
df.groupby(["letter", "one"], as index=False).sum
()
```

	letter	one	two
0	а	1	4
1	b	1	4
2	С	1	2

		two
letter	one	
a	1	4
b	1	4
С	1	2

Lee 704 Pandas

Missing Values: numpy.nan and pandas.NA

```
>>> import numpy as np
>>> df = pd.DataFrame(np.random.randn(5,3),
        index=['a','c','e','f','h'],
        columns=["one","two","three"])
>>> df.loc['a',"two"] = np.nan
>>> df
                two three
       one
              NaN 2.291062
a 0.405857
c - 0.491699 1.178464 - 1.951439
e 2.649431 -1.938993 -0.827796
f 1.235505 -0.726769 -0.146558
h 2.837826 0.505976 -0.489474
```



Lee 704 Pandas 27 ROBERT H. SMITH SCHOOL OF BUSINES.

Detect Missing Values

- pandas.isnull() Detect missing values for an array-like object
- pandas.Series.isnull() Detect missing values in a Series
- pandas.DataFrame.isnull() Detect missing values in a DataFrame

```
>>> df.isnull()
                three
            t.wo
    one
  False
          True
                False
а
         False
  False
                False
  False
         False
                False
f
  False
         False
                False
  False
         False
                False
```



Lee 704 Pandas

28 ROBERT H. SMITH

Fill Missing Values

```
>>> df.fillna(0)

one two three

a 0.405857 0.000000 2.291062

c -0.491699 1.178464 -1.951439

e 2.649431 -1.938993 -0.827796

f 1.235505 -0.726769 -0.146558

h 2.837826 0.505976 -0.489474
```



Lee 704 Pandas

Query Data

Use the query(expression) function to embed boolean expressions on columns within quotes:

```
>>> df.query("one>0")
                 t.wo
                        three
       one
  0.405857 NaN 2.291062
а
e 2.649431 -1.938993 -0.827796
f 1.235505 -0.726769 -0.146558
h 2.837826 0.505976 -0.489474
>>> df.query("one>0 & two>0")
                        three
                 two
       one
h 2.837826 0.505976 -0.489474
```



Lee 704 Pandas 30 ROBERT H. SMITH

Query Data

You can apply any function to the columns:

You can apply any function to element-wise:

```
>>> df.applymap(np.sqrt)
                 t.wo
                         three
       one
  0.637069
           NaN 1.513626
       NaN 1.085571
                           NaN
e 1.627707
                 NaN
                           NaN
f 1.111533
                 NaN
                           NaN
h 1.684585 0.711320
                           NaN
```



Lee 704 Pandas

Pandas - Visualization

http://pandas.pydata.org/pandas-docs/stable/ visualization.html

- Use the standard convention for referencing the matplotlib API
- Provide the basics in pandas to easily create decent looking plots



Lee 704 Pandas 32 ROBERT H. SMITH