

PANDAS

A PYTHON DATA ANALYSIS LIBRARY

- Built on top of numpy to make data analysis easier
- Automatic data alignment based on labels or indices
- Data aggregation, transformation and grouping
- Intuitive merging and joining of datasets
- Hierarchical labeling
- Reading and Writing of CSV, Excel and others

PANDAS.SERIES

- For storing indexed 1D data

creation from numpy array with list as index

```
s = pd.Series(np.arange(5), index=['a', 'b', 'c', 'd', 'e'])  
print(s)
```

```
a      0  
b      1  
c      2  
d      3  
e      4  
dtype: int64
```

INDEX IS CREATED IF NOT SET

```
pd.Series(randn(5))
```

```
0    -0.683954  
1    -0.653830  
2     0.712992  
3     0.333370  
4    -0.769677  
dtype: float64
```

SERIES IS LIKE AN ARRAY

```
s[0]  
print("\n")  
s[s > s.median()]  
print("\n")  
s[[3,2,1]]
```

```
0  
  
d      3  
e      4  
dtype: int64  
  
d      3  
c      2  
b      1  
dtype: int64
```

SERIES IS LIKE A DICTIONARY

```
s['a']  
s['e'] = 6  
s  
'e' in s  
'f' in s
```

```
0  
>>> a      0  
b      1  
c      2  
d      3  
e      6  
dtype: int64  
True  
False
```

OPERATIONS ON SERIES

```
s+s  
s**2  
np.exp(s)
```

```
a      0  
b      2  
c      4  
d      6  
e     12  
dtype: int64  
a      0  
b      1  
c      4  
d      9  
e     36  
dtype: int64  
a      1.000000  
b      2.718282  
c      7.389056  
d     20.085537  
e    403.428793  
dtype: float64
```

PANDAS.DATFRAME

A 2D labeled data structure with columns of potentially different types.

Like Series, DataFrame accepts many different kinds of input:

- Dict of 1D ndarrays, lists, dicts, or Series
- 2-D numpy.ndarray
- Structured or record ndarray
- A Series
- Another DataFrame

FROM DICTIONARY

```
d = {'one' : pd.Series([1., 2., 3.], index=['a', 'b', 'c']),  
     'two' : pd.Series([1., 2., 3., 4.], index=['a', 'b', 'c', 'd'])  
}  
df = pd.DataFrame(d)
```

df

	one	two
a	1	1
b	2	2
c	3	3
d	NaN	4

FROM OTHER DATAFRAME

```
pd.DataFrame(df, index=['d', 'b', 'a'])
```

	one	two
d	NaN	4
b	2	2
a	1	1

```
pd.DataFrame(d, index=['d', 'b', 'a'], columns=['two', 'three'])
```

	two	three
d	4	NaN
b	2	NaN
a	1	NaN

COMPLEX CASES

```
df2 = pd.DataFrame({'A': 1.,  
                    'B': pd.Timestamp('20130102'),  
                    'C': pd.Series(1,index=list(range(4)),  
                                   dtype='float32'),  
                    'D': np.array([3] * 4,dtype='int32'),  
                    'E': 'foo' })
```

df2

	A	B	C	D	E
0	1	2013-01-02	1	3	foo
1	1	2013-01-02	1	3	foo
2	1	2013-01-02	1	3	foo
3	1	2013-01-02	1	3	foo

```
df2 = pd.DataFrame({'A': 1.,  
                    'B': pd.Timestamp('20130102'),  
                    'C': pd.Series(1,index=list(range(4)),  
                                   dtype='float32'),  
                    'D': np.array([3] * 4,dtype='int32'),  
                    'E': 'foo' })
```

```
df2.dtypes
```

```
A          float64  
B    datetime64[ns]  
C          float32  
D          int32  
E          object  
dtype: object
```

TIME SERIES

```
# Date range
dates = pd.date_range('20130101', periods=6)
# Dataframes
df = pd.DataFrame(np.random.randn(6, 4), index=dates, columns=list('ABCD'))
```

df

	A	B	C	D
2013-01-01	1.295060	-1.892445	1.325433	-2.208407
2013-01-02	-0.188383	0.957533	1.024778	0.827907
2013-01-03	1.005206	-0.734869	0.685699	-0.092167
2013-01-04	-0.882230	-0.256595	1.627533	-0.509552
2013-01-05	-0.193930	2.171011	-1.064996	1.736236
2013-01-06	-1.995633	-0.536571	-0.043122	0.773973

INSPECTION

```
df.head()
```

	A	B	C	D
2013-01-01	1.295060	-1.892445	1.325433	-2.208407
2013-01-02	-0.188383	0.957533	1.024778	0.827907
2013-01-03	1.005206	-0.734869	0.685699	-0.092167
2013-01-04	-0.882230	-0.256595	1.627533	-0.509552
2013-01-05	-0.193930	2.171011	-1.064996	1.736236

```
df.tail(3)
```

	A	B	C	D
2013-01-04	-0.882230	-0.256595	1.627533	-0.509552
2013-01-05	-0.193930	2.171011	-1.064996	1.736236
2013-01-06	-1.995633	-0.536571	-0.043122	0.773973

COLUMNS AND VALUES

```
df.columns, df.values
```

```
(Index(['A', 'B', 'C', 'D'], dtype='object'),  
array([[ 1.29505997, -1.89244519,  1.32543257, -2.20840719],  
       [-0.1883832 ,  0.95753289,  1.02477782,  0.82790719],  
       [ 1.00520553, -0.73486871,  0.68569858, -0.09216724],  
       [-0.88222997, -0.25659454,  1.62753267, -0.50955237],  
       [-0.19392973,  2.17101115, -1.06499636,  1.73623578],  
       [-1.9956328 , -0.53657104, -0.0431219 ,  0.77397318]]))
```

DESCRIBE A DATAFRAME

```
df.describe()
```

	A	B	C	D
count	6.000000	6.000000	6.000000	6.000000
mean	-0.159985	-0.048656	0.592554	0.087998
std	1.213921	1.420640	0.995524	1.370964
min	-1.995633	-1.892445	-1.064996	-2.208407
25%	-0.710155	-0.685294	0.139083	-0.405206
50%	-0.191156	-0.396583	0.855238	0.340903
75%	0.706808	0.654001	1.250269	0.814424
max	1.295060	2.171011	1.627533	1.736236

DATAFRAME SLICING OVERVIEW

Operation	Syntax	Result
Select column	<code>df[col]</code>	Series
Select row by label	<code>df.loc[label]</code>	Series
Select row by integer location	<code>df.iloc[loc]</code>	Series
Slice rows	<code>df[5:10]</code>	DataFrame
Select rows by boolean vector	<code>df[bool_vec]</code>	DataFrame

BY COLUMN OR ROW SLICE

```
df['A']
```

```
2013-01-01    1.295060
2013-01-02   -0.188383
2013-01-03    1.005206
2013-01-04   -0.882230
2013-01-05   -0.193930
2013-01-06   -1.995633
Freq: D, Name: A, dtype: float64
```

```
df[0:3]
```

	A	B	C	D
2013-01-01	1.295060	-1.892445	1.325433	-2.208407
2013-01-02	-0.188383	0.957533	1.024778	0.827907
2013-01-03	1.005206	-0.734869	0.685699	-0.092167

BY INDEX

```
df['20130102':'20130104']
```

	A	B	C	D
2013-01-02	-0.188383	0.957533	1.024778	0.827907
2013-01-03	1.005206	-0.734869	0.685699	-0.092167
2013-01-04	-0.882230	-0.256595	1.627533	-0.509552

```
from datetime import date  
df[date(2013,1,2):date(2013,1,4)]
```

	A	B	C	D
2013-01-02	-0.188383	0.957533	1.024778	0.827907
2013-01-03	1.005206	-0.734869	0.685699	-0.092167
2013-01-04	-0.882230	-0.256595	1.627533	-0.509552

BY INTEGER LOCATION

```
df.iloc[[4, 2]]
```

	A	B	C	D
2013-01-05	-0.193930	2.171011	-1.064996	1.736236
2013-01-03	1.005206	-0.734869	0.685699	-0.092167

GROUPING

```
gp = pd.DataFrame({'A' : ['foo', 'bar', 'foo', 'bar',  
                          'foo', 'bar', 'foo', 'foo'],  
                  'B' : ['one', 'one', 'two', 'three',  
                          'two', 'two', 'one', 'three'],  
                  'C' : np.random.randn(8),  
                  'D' : np.random.randn(8)})
```

gp

	A	B	C	D
0	foo	one	0.771749	-1.173628
1	bar	one	0.005310	0.168708
2	foo	two	-1.281703	-0.365025
3	bar	three	0.296035	-0.569597
4	foo	two	0.640401	0.015837
5	bar	two	-1.178718	-0.531231
6	foo	one	-0.875578	0.047414
7	foo	three	0.269448	1.214574

```
gp.groupby('A').sum()
```

	C	D
A		
bar	-0.877373	-0.932120
foo	-0.475683	-0.260827

```
gp.groupby(['A', 'B']).mean()
```

		C	D
A	B		
bar	one	0.005310	0.168708
	three	0.296035	-0.569597
	two	-1.178718	-0.531231
foo	one	-0.051915	-0.563107
	three	0.269448	1.214574
	two	-0.320651	-0.174594

MERGING

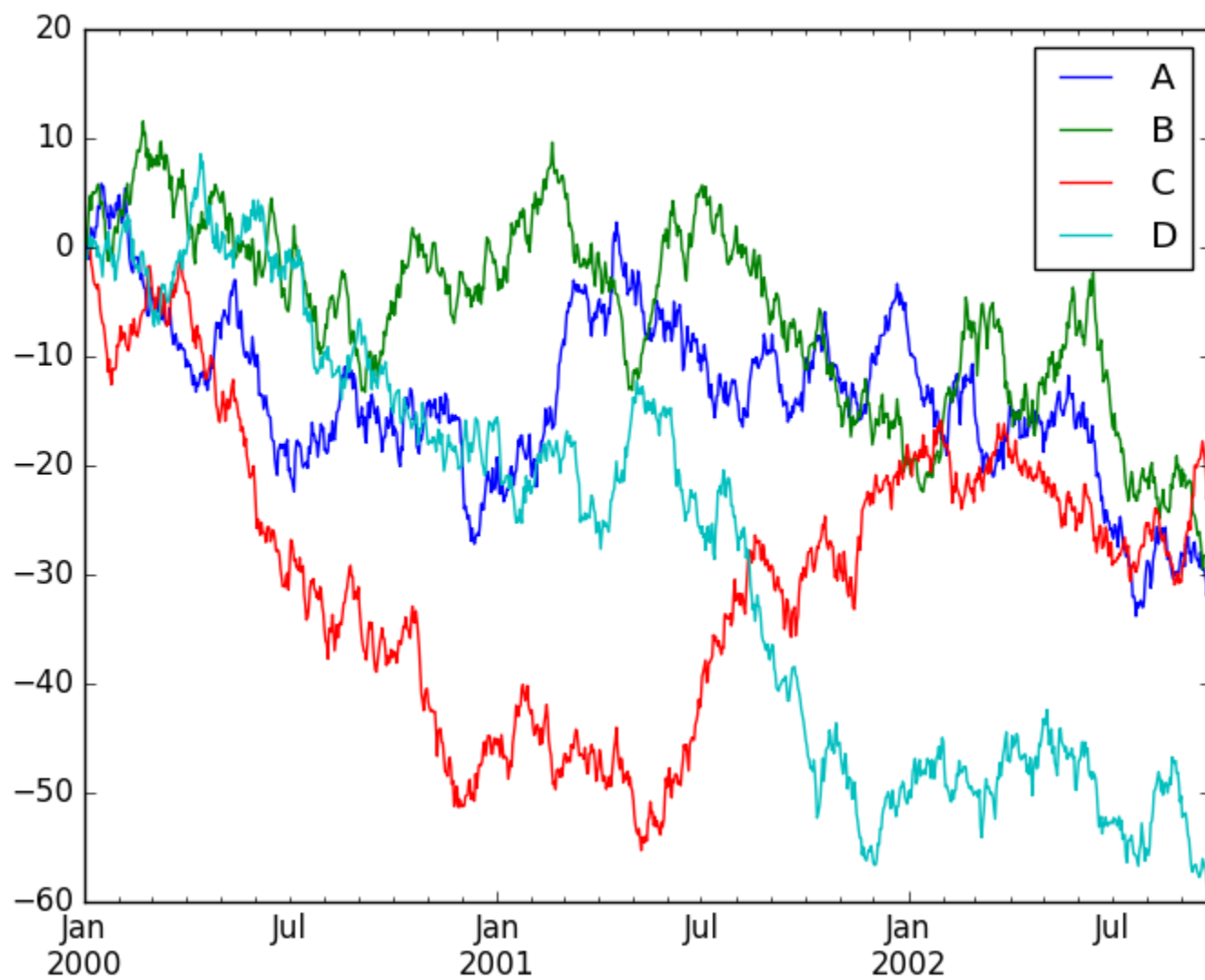
```
left = pd.DataFrame({'key': ['one', 'two'], 'lval': [1, 2]})  
right = pd.DataFrame({'key': ['two', 'one'], 'rval': [4, 5]})  
pd.merge(left, right, on='key')
```

	key	lval	rval
0	one	1	5
1	two	2	4

PLOTTING

Pandas has built-in functions for common plot types

```
import matplotlib.pyplot as plt
df = pd.DataFrame(randn(1000, 4),
                  index=pd.date_range('1/1/2000', periods=1000),
                  columns=list('ABCD'))
df = df.cumsum()
ax = df.plot()
```

WORKING WITH A DATASET

Let's try working with the [Movielens](#) 100k dataset

- 1000 Users
- 100,000 Ratings
- 1700 Movies

Extract the ml-100k.zip to a folder `ml - 100k` in the same directory as the `lecture7.py`

READING THE DATA

```
# pass in column names for each CSV
u_cols = ['user_id', 'age', 'sex', 'occupation', 'zip_code']
users = pd.read_csv('ml-100k/u.user', sep='|', names=u_cols,
                    encoding="latin-1")

r_cols = ['user_id', 'movie_id', 'rating', 'unix_timestamp']
ratings = pd.read_csv('ml-100k/u.data', sep='\t', names=r_cols,
                      encoding="latin-1")

# the movies file contains columns indicating the movie's genres
# let's only load the first five columns of the file with usecols
m_cols = ['movie_id', 'title', 'release_date',
          'video_release_date', 'imdb_url']
movies = pd.read_csv('ml-100k/u.item', sep='|',
                     names=m_cols, usecols=range(5),
                     encoding='latin-1')

# create one merged DataFrame
movie_ratings = pd.merge(movies, ratings)
lens = pd.merge(movie_ratings, users)
```

WHAT DID WE READ?

```
lens.head(3)
```

```

  movie_id      title  release_date  video_release_date  \
0         1  Toy Story (1995)    01-Jan-1995           NaN
1         4  Get Shorty (1995)    01-Jan-1995           NaN
2         5   Copycat (1995)    01-Jan-1995           NaN

                                imdb_url  user_id  rating
g  \
0  http://us.imdb.com/M/title-exact?Toy%20Story%2...    308
4
1  http://us.imdb.com/M/title-exact?Get%20Shorty%...    308
5
2  http://us.imdb.com/M/title-exact?Copycat%20(1995)    308
4

  unix_timestamp  age  sex  occupation  zip_code
0      887736532   60   M    retired    95076
1      887737890   60   M    retired    95076
2      887739608   60   M    retired    95076
```

WHAT ARE THE 10 MOST RATED MOVIES?

```
most Rated = lens.groupby('title').size().sort_values(ascending=False)[:10]
print(most Rated)
```

```
title
Star Wars (1977)      583
Contact (1997)        509
Fargo (1996)          508
Return of the Jedi (1983) 507
Liar Liar (1997)      485
English Patient, The (1996) 481
Scream (1996)         478
Toy Story (1995)      452
Air Force One (1997)  431
Independence Day (ID4) (1996) 429
dtype: int64
```

WHICH MOVIES ARE MOST HIGHLY RATED?

The `agg` function can take multiple functions that are applied to a column

```
movie_stats = lens.groupby('title').agg({'rating': [np.size, np.mean]})  
movie_stats.head()
```

	rating size	mean
title		
'Til There Was You (1997)	9	2.333333
1-900 (1994)	5	2.600000
101 Dalmatians (1996)	109	2.908257
12 Angry Men (1957)	125	4.344000
187 (1997)	41	3.024390

WHICH MOVIES ARE MOST HIGHLY RATED?

Sort them by mean rating

```
movie_stats.sort([('rating', 'mean')], ascending=False).head()
```

title	rating size	mean
They Made Me a Criminal (1939)	1	5
Marlene Dietrich: Shadow and Light (1996)	1	5
Saint of Fort Washington, The (1993)	2	5
Someone Else's America (1995)	1	5
Star Kid (1997)	3	5

WHICH MOVIES ARE MOST HIGHLY RATED?

Lets only look at movies rated at least 100 times

```
atleast_100 = movie_stats['rating'].size >= 100  
movie_stats[atleast_100].sort([('rating', 'mean')], ascending=False)  
.head()
```

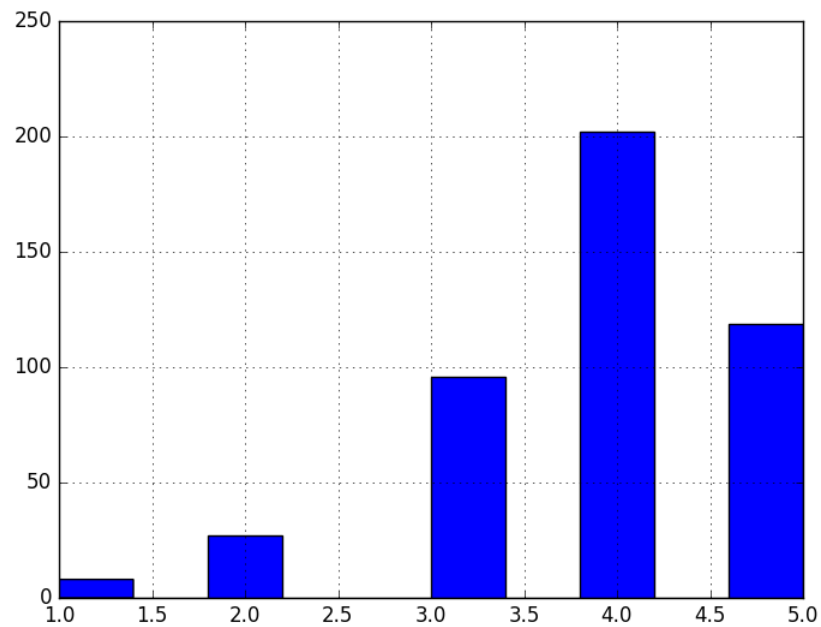
```
'org_babel_python_eoe'
```


EXERCISE

```
### Exercise ###  
### Try to plot the ratings distribution of a movie of your choice.  
### you can use the hist() function to produce a histogram
```

SOLUTION

```
toy_story = lens[lens.title=='Toy Story (1995)']  
plt.figure()  
ax = toy_story.rating.hist()  
plt.savefig('hist.png')  
'hist.png'
```

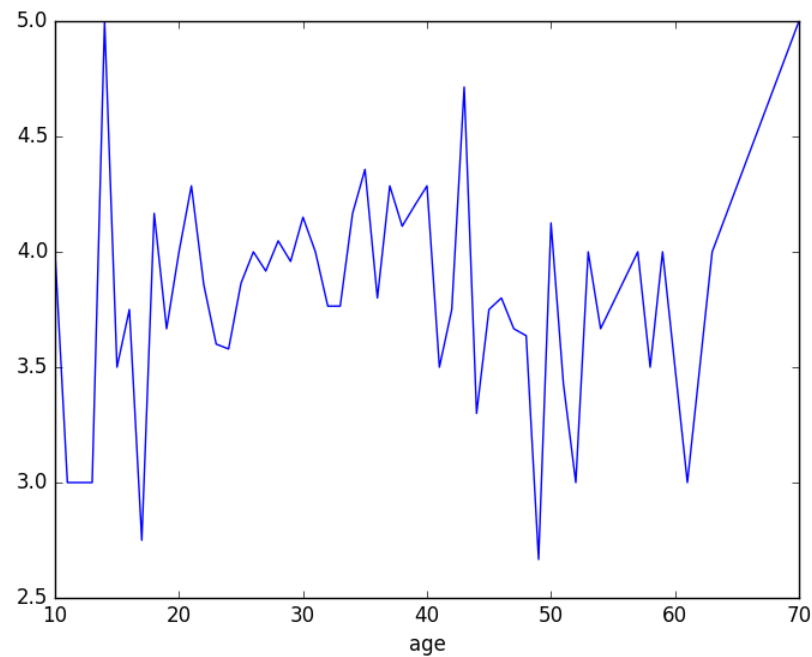


EXERCISE 2

```
### Exercise ###  
### plot the mean rating by age of user
```

SOLUTION

```
age_grouped = toy_story.groupby('age').mean()  
plt.figure()  
ax = age_grouped['rating'].plot()  
plt.savefig('age-ratings.png')  
'age-ratings.png'
```



ADDITIONAL RESOURCES

- [Pandas website](#) - The documentation is very thorough and full of examples
- [List of pandas tutorials](#)
- [using pandas on the movielens dataset](#) (blogpost from which I took some examples)