

# PYTHON FOR FINANCE

---

*Analyze Big Financial Data*



# Introduction

- How information influences the movement of stock
- Building software to analyze & visualize the relationships
- Details about how stock exchange works
- Learn how to build machine learning algorithms

# Three parts to the course

- Manipulating Financial Data in Python
  - How read data into Python and manipulate using power statistical algorithms
- Computational Investing
  - Algorithms, methods & models used by hedge funds and IBs to manipulate and work with financial data
- Learning Algorithms for Trading
  - We pull everything together, the data and use it with machine learning algorithm like Q-Learning and Random Forest.

# Goal

- After completion of the 'Series' of workshops:
  - You'll be equipped to join trading system development team.
    - Hedge fund or Investment Bank
  - Begin equity trading, and earn some money !
  - Build up your own portfolio of stocks and other financial instruments.

# Textbooks

- Python for Finance by Hilpish
- What Hedge Funds Really Do by Philip J. Romero
- Machine Learning by Mitchell



# Python for Financial Applications:

- Set up quick prototype algorithms
- Computational speed
- Features:
  - Strong scientific libraries
  - Strongly maintained
  - Fast
    - If you stick to metrics notation, because lower levels are returned in C.
- Alternatives: R and MATLAB

# Prerequisites

- You need to have strong programming skills
- Basic knowledge of finance
- Statistics/Econometrics knowledge

Prepare for the challenge !

# READING & PLOTTING OF STOCK DATA

---



# Lesson outline

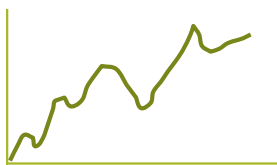
- In this lesson you will learn how to read data, select subsets of it and generate useful plots, using [pandas](#) and [matplotlib](#).
- Read stock data from CSV files:
  - [pandas.DataFrame](#)
  - [pandas.read\\_csv](#)
- Select desired rows and columns:
  - [Indexing and Slicing Data](#)
  - Gotchas: [Label-based slicing conventions](#)
- Visualize data by generating plots:
  - [Plotting](#)
  - [pandas.DataFrame.plot](#)
  - [matplotlib.pyplot.plot](#)

# What CSV file looks

```
Date,Open,High,Low,Close,Adj Close,  
2006-11-30,5486.000000,5684.399902,5411.000000,5669.899902,5669.899902,0  
2006-12-31,5672.700195,5829.299805,5499.000000,5773.399902,5773.399902,0  
2007-01-31,5781.299805,6052.100098,5781.299805,5832.500000,5832.500000,0  
2007-02-28,5840.600098,5998.899902,5641.100098,5995.000000,5995.000000,0  
2007-03-31,5994.000000,6265.200195,5915.799805,6166.000000,6166.000000,0  
2007-04-30,6163.200195,6384.100098,6134.200195,6313.500000,6313.500000,0  
2007-05-31,6321.000000,6409.200195,6170.200195,6274.899902,6274.899902,0  
2007-06-30,6282.000000,6436.700195,6044.299805,6144.200195,6144.200195,0  
2007-07-31,6138.600098,6247.200195,5483.299805,6247.200195,6247.200195,0  
2007-08-31,6253.700195,6594.399902,6146.500000,6567.799805,6567.799805,0  
2007-09-30,6567.600098,6800.200195,6542.899902,6754.100098,6754.100098,0  
2007-10-31,6778.399902,6851.500000,6312.600098,6533.100098,6533.100098,0  
2007-11-30,6535.200195,6684.399902,6105.600098,6339.799805,6339.799805,0  
2007-12-31,6336.600098,6385.700195,5186.799805,5650.299805,5650.299805,0  
2008-01-31,5672.600098,6022.000000,5490.399902,5572.100098,5572.100098,0  
2008-02-29,5534.299805,5535.200195,5039.600098,5355.700195,5355.700195,0  
2008-03-31,5363.399902,5672.700195,5301.000000,5595.399902,5595.399902,0  
2008-04-30,5592.600098,5980.799805,5538.700195,5654.700195,5654.700195,0  
2008-05-31,5671.899902,5681.500000,5144.700195,5215.299805,5215.299805,0
```

Header

Data



# Quiz:

- Which fields you would expect to see in a CSV file of stock data?

☐ # of employees

☐ Date/time

☐ Company name

☐ Price of the stock

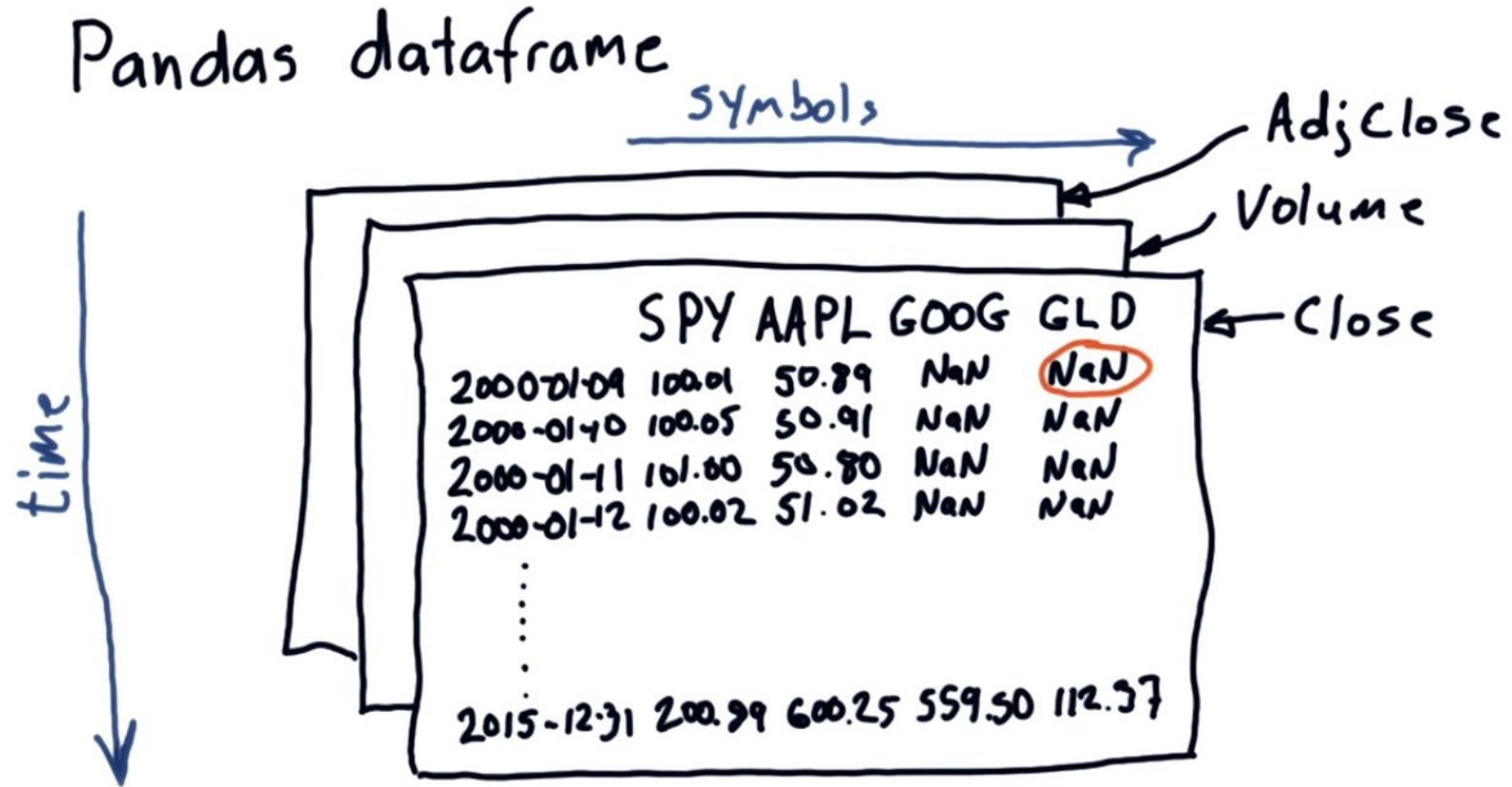
☐ Company's hometown

# Yahoo Database

A handwritten diagram of a database table with three rows. The columns are labeled: Date, Open, High, Low, Close, Volume, and Adj Close. The first row is dated 2012-01-11 and has values 46.23, 46.78, 46.20, 46.73, 1955400, and 46.73. The second row is dated 2012-08-09 and has values 45.56, 45.93, 45.47, 45.37, 1963000, and 45.57. The third row is dated 2000-02-01 and has values 24.37, 25.15, 24.44, 25.00, 413200, and 5.36. A vertical arrow on the left points upwards from 'older' to 'newer'. A red bracket under the 'Close' values of the first and third rows is labeled 'same', and another red bracket under the 'Volume' and 'Adj Close' values of the first and third rows is labeled 'differ'.

	<u>Date</u>	<u>Open</u>	<u>High</u>	<u>Low</u>	<u>Close</u>	<u>Volume</u>	<u>Adj Close</u>
newer ↑	<u>2012-01-11</u>	46.23	46.78	46.20	46.73	1955400	46.73
	⋮	⋮	⋮	⋮	⋮	⋮	⋮
	2012-08-09	45.56	45.93	45.47	45.37	1963000	45.57
	⋮	⋮	⋮	⋮	⋮	⋮	⋮
older	<u>2000-02-01</u>	24.37	25.15	24.44	25.00	413200	5.36

# Pandas DataFrame



# Read CSV file

```
1 import pandas as pd
2
3 def test_run():
4     df=pd.read_csv("C:\Users\Fahad\Desktop\Data\AAPL.csv")
5     print df
6     print df.head()
7     
8     print df[0:5]
9
10 if __name__=="__main__":
11     test_run()
```

# Quiz: Print last 5 rows of the data frame

```
1 import pandas as pd
2
3
4 def test_run():
5     """Function called by Test Run."""
6     df = pd.read_csv("data/AAPL.csv")
7     # TODO: Print last 5 rows of the data frame
8
9
10 if __name__ == "__main__":
11     test_run()
12
```

# Slicing

```
1 import pandas as pd
2
3
4 def test_run():
5     df = pd.read_csv("data/AAPL.csv")
6     print df[10:21]#rows between index 10 and 20
7
8
9 if __name__ == "__main__":
10     test_run()
```



# Computing max closing price

```
1 import pandas as pd
2
3 def get_max_close(symbol):
4     """Return the maximum closing value for stock indicated by symbol.
5
6     Note: Data for a stock is stored in file: data/<symbol>.csv
7     """
8     df = pd.read_csv("data/{}.csv".format(symbol)) # read in data
9     return df['Close'].max() # compute and return max
10
11
12 def test_run():
13     """Function called by Test Run."""
14     for symbol in ['AAPL', 'IBM']:
15         print "Max close"
16         print symbol, get_max_close(symbol)
17
18
19
20 if __name__ == "__main__": # if run standalone
21     test_run()
```

# Quiz: Compute mean volume

```
1  """Compute mean volume"""
2
3  import pandas as pd
4
5  def get_mean_volume(symbol):
6      """Return the mean volume for stock indicated by symbol.
7
8      Note: Data for a stock is stored in file: data/<symbol>.csv
9      """
10     df = pd.read_csv("data/{}.csv".format(symbol)) # read in data
11     # TODO: Compute and return the mean volume for this stock
12
13
14  def test_run():
15      """Function called by Test Run."""
16      for symbol in ['AAPL', 'IBM']:
17          print "Mean Volume"
18          print symbol, get_mean_volume(symbol)
19
20
21  if __name__ == "__main__":
22      test_run()
23
```

# Plotting stock price data

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 def test_run():
5     df = pd.read_csv("data/AAPL.csv")
6     print df['Adj Close']
7     df['Adj Close'].plot()
8     plt.show() # must be called to show plots
9
10
11 if __name__ == "__main__":
12     test_run()
```

# Quiz: Plot High prices for IBM

```
1  """Plot High prices for IBM"""
2
3  import pandas as pd
4  import matplotlib.pyplot as plt
5
6  def test_run():
7      df = pd.read_csv("data/IBM.csv")
8      # TODO: Your code here
9      plt.show() # must be called to show plots
10
11
12  if __name__ == "__main__":
13      test_run()
14
```

# Plot two columns

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 def test_run():
5     df = pd.read_csv("data/AAPL.csv")
6     df[['Close', 'Adj Close']].plot()
7     plt.show() # must be called to show plots
8
9
10 if __name__ == "__main__":
11     test_run()
```

# WORKING WITH MULTIPLE STOCKS

---

# Lesson Outline

- Here's an overview of what you'll learn to do in this lesson. Documentation links are for reference.
- **Read in multiple stocks:**
  - Create an empty [`pandas.DataFrame`](#) with dates as index: [`pandas.date\_range`](#)
  - Drop missing date rows: [`pandas.DataFrame.dropna`](#)
  - Incrementally join data for each stock: [`pandas.DataFrame.join`](#)
- **Manipulate stock data:**
  - [`Index and select data`](#) by row (dates) and column (symbols)
  - Plot multiple stocks at once (still using [`pandas.DataFrame.plot`](#))
  - Carry out arithmetic operations across stocks

# Problems to solve

Pandas dataframe

	<u>Symbols</u>	SPY	IBM	GOOG	GLD
2010-1-22		100.10	99.09	500	20.10
2010-1-25		100.15	100.02	560	20.09
...		...	...	...	...
2012-12-31		200.15	150.01	100.09	30.00

Problems to solve

- Date ranges
- Multiple stocks
- Align dates
- Proper date order



## Quiz: NYSE trading days

Quiz: How many days were US  
stocks traded at NYSE in 2014?

✱ 365

✱ 260

✱ 252

# Building a DataFrame

Building a dataframe

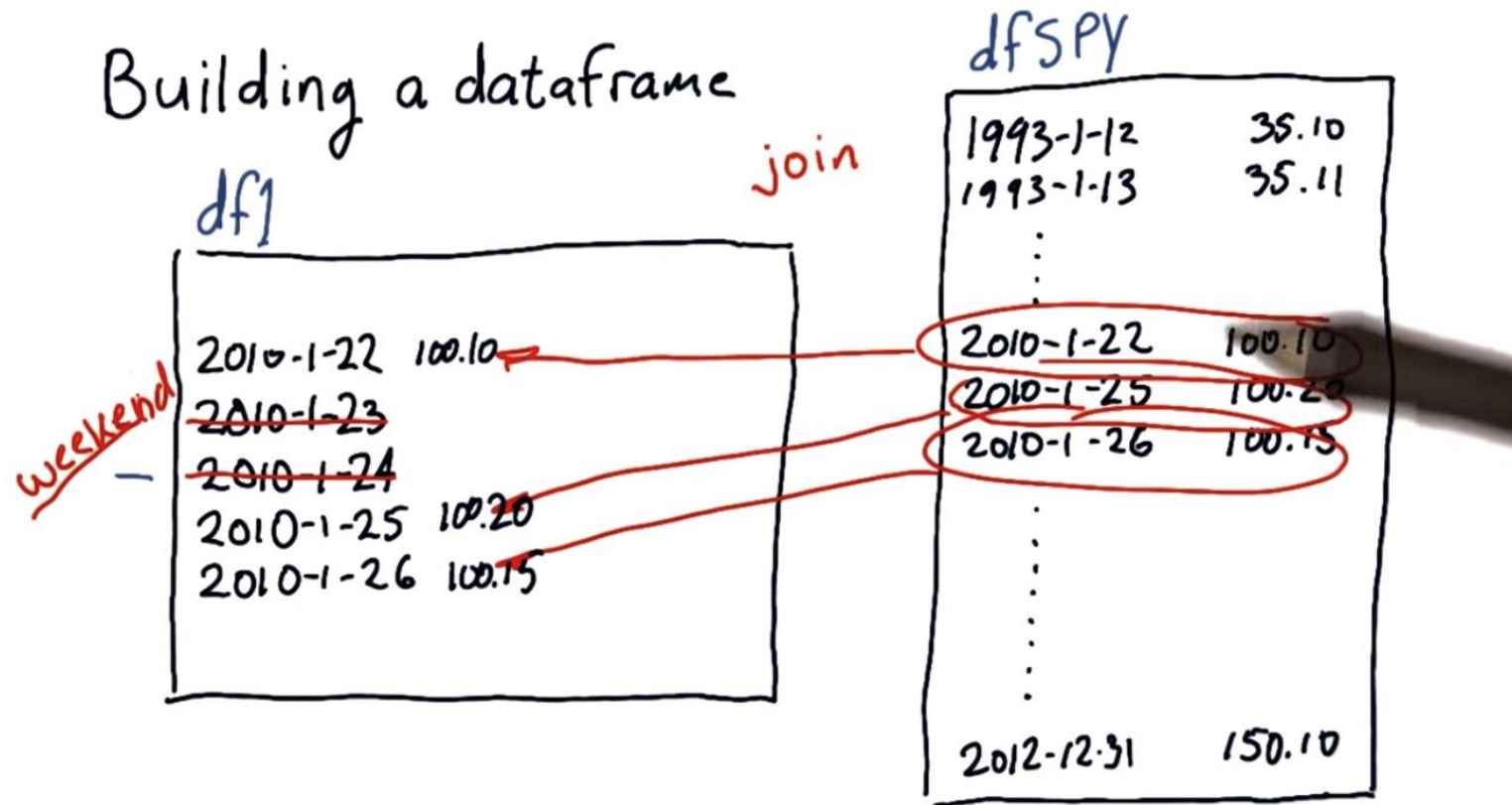
df1

	2010-1-22
-	2010-1-23
-	2010-1-24
	2010-1-25
	2010-1-26

dfSPY

1993-1-12	35.10
1993-1-13	35.11
...	
2010-1-22	100.10
2010-1-25	100.20
2010-1-26	100.15
...	
2012-12-31	150.10

# Building a DataFrame



# Building a DataFrame: 'JOIN'

Building a dataframe

df1

	SPY	IBM
2010-1-22	100.10	120.21
2010-1-25	100.20	120.22
2010-1-26	100.15	120.10

join

dfIBM

1993-1-12	80.19
1993-1-13	80.20
...	...
2010-1-22	120.21
2010-1-25	120.22
2010-1-26	120.10
...	...
2012-12-31	150.29

# Create an empty data frame

```
1 import pandas as pd
2 def test_run():
3     start_date='2016-10-26'
4     end_date='2016-11-09'
5     dates=pd.date_range(start_date,end_date)
6     print dates[0]
7     df1=pd.DataFrame(index=dates)
8     print df1
9 if __name__ == "__main__":
10     test_run()
11
12
13
```

# Join SPY data

```
4 def test_run():
5     #Define date range
6     start_date='2010-01-22'
7     end_date='2010-01-26'
8     dates=pd.date_range(start_date,end_date)
9
10    #Create an empty dataframe
11    df1=pd.DataFrame(index=dates)
12
13    #Read SPY data into temporary dataframe
14    dfSPY = pd.read_csv("data/SPY.csv",index_col="Date",
15                        parse_dates=True,usecols=['Date','Adj Close'],
16                        na_values=['nan'])
17
18    #Join the two dataframes using DataFrame.join()
19    df1=df1.join(dfSPY)
20
21    # Drop NaN Values
22    df1 = df1.dropna()
23    print df1
```

# Read in more stocks

```
12
13     #Read SPY data into temporary dataframe
14     dfSPY = pd.read_csv("data/SPY.csv",index_col="Date",
15                        parse_dates=True,usecols=['Date','Adj Close'],
16                        na_values=['nan'])
17
18     #Rename 'Adj Close' column to 'SPY' to prevent clash
19     dfSPY = dfSPY.rename(columns={'Adj Close':'SPY'})
20
21     #Join the two dataframes using DataFrame.join(), with how='inner'
22     df1=df1.join(dfSPY,how='inner')
23
24     #Read in more stocks
25     symbols = ['GOOG','IBM','GLD']
26     for symbol in symbols:
27         df_temp=pd.read_csv("data/{}.csv".format(symbol), index_col='Date',
28                            parse_dates=True,usecols=['Date','Adj Close']
29                            ,na_values=['nan'])
30
31         # rename to prevent clash
32         df_temp = df_temp.rename(columns={'Adj Close': symbol}) |
33         df1=df1.join(df_temp) #use default how='left'
```

# Quiz: Utility functions for reading data

```
5
6 ▾ def symbol_to_path(symbol, base_dir="data"):
7     """Return CSV file path given ticker symbol."""
8     return os.path.join(base_dir, "{}.csv".format(str(symbol)))
9
10
11 ▾ def get_data(symbols, dates):
12     """Read stock data (adjusted close) for given symbols from CSV files."""
13     df = pd.DataFrame(index=dates)
14 ▾     if 'SPY' not in symbols: # add SPY for reference, if absent
15         symbols.insert(0, 'SPY')
16
17 ▾     for symbol in symbols:
18         # TODO: Read and join data for each symbol
19
20     return df
21
22
23 ▾ def test_run():
24     # Define a date range
25     dates = pd.date_range('2010-01-22', '2010-01-26')
26
27     # Choose stock symbols to read
28     symbols = ['GOOG', 'IBM', 'GLD']
29
30     # Get stock data
31     df = get_data(symbols, dates)
32     print df
33
34
```



# Solution

```
4 import pandas as pd
5
6 def symbol_to_path(symbol, base_dir="data"):
7     """Return CSV file path given ticker symbol."""
8     return os.path.join(base_dir, "{}.csv".format(str(symbol)))
9
10
11 def get_data(symbols, dates):
12     """Read stock data (adjusted close) for given symbols from CSV files."""
13     df = pd.DataFrame(index=dates)
14     if 'SPY' not in symbols: # add SPY for reference, if absent
15         symbols.insert(0, 'SPY')
16
17     for symbol in symbols:
18         df_temp = pd.read_csv(symbol_to_path(symbol), index_col='Date',
19                               parse_dates=True, usecols=['Date', 'Adj Close'], na_values=['nan'])
20         df_temp = df_temp.rename(columns = {'Adj Close': symbol})
21         df = df.join(df_temp)
22         if symbol == 'SPY': #drop dates SPY did not trade
23             df = df.dropna(subset=["SPY"])
24
25
```

# Slicing DataFrame

Slicing dataframes

`df2` = `df1[sliced, ['GOOG', 'GLD']]`

	GOOG	GLD
2010-2-13	~	~
2010-2-14	~	~
2010-2-15	~	~

	SPY	IBM	GOOG	GLD
2010-1-22	~	~	~	~
2010-1-25	~	~	~	~
...				
2010-2-13	~	~	~	~
2010-2-14	~	~	~	~
2010-2-15	~	~	~	~
...				
2012-12-31	~	~	~	~

The diagram illustrates the slicing operation. A red box in the `df1` table highlights the rows for dates 2010-2-13, 2010-2-14, and 2010-2-15, and the columns for GOOG and GLD. A red arrow points from this box to the `df2` table, which contains only the sliced data.

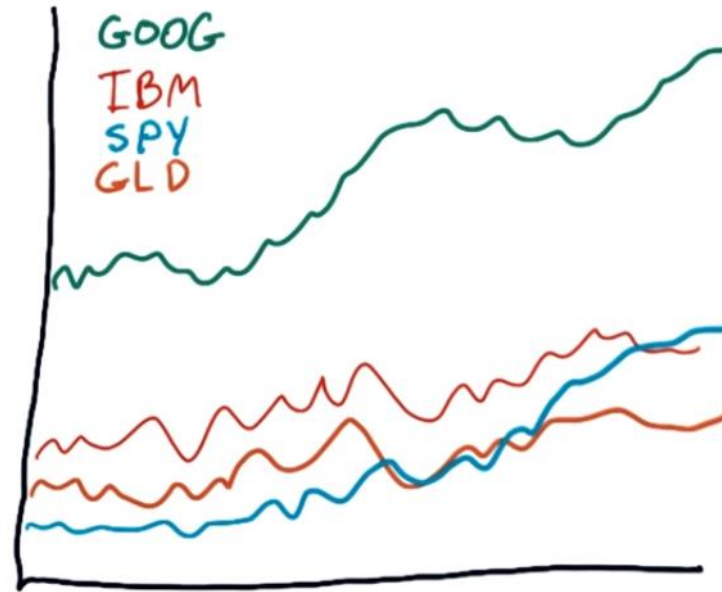
# More Slicing

```
31
32     # Get stock data
33     df = get_data(symbols, dates)
34     print df
35     # More slicing
36     print
37     print 'More slicing'
38     print df.ix['2012-10-30': '2012-11-04']
39     print
40     print df['IBM']    #print one label
41     print
42     print df[['IBM', 'AAPL']] #print 2 columns
43     #Slice by row & column:
44     print 'Slice by row & column:'
45     print df.ix['2012-10-28':'2012-11-04',['IBM', 'AAPL']]
46     #normalize
47     df_n = df / df.ix[0]
48     print df_n
49
```

# Problems with plotting

Plotting a dataframe  
df.plot()

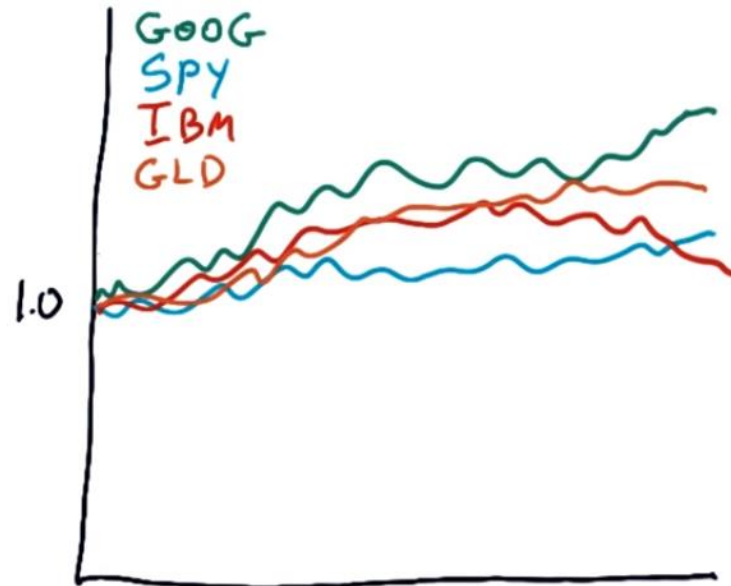
	SPY	IBM	GOOG	GLD
2010-1-22	~	~	~	~
~~~~	~	~	~	~
~~~~	~	~	~	~
~~~~	~	~	~	~
2012-12-31	~	~	~	~



# Normalization

Plotting a dataframe  
df1.plot()

	SPY	IBM	GOOG	GLD
2010-1-22	~	~	~	~
~~~~~	~	~	~	~
~~~~~	~	~	~	~
~~~~~	~	~	~	~
~~~~~	~	~	~	~
2012-12-31	~	~	~	~



## Quiz: How to plot on "equal footing"?

Quiz: What is the best way  
to normalize price data  
so that all prices start at 1.0?

★ `for date in df1.index:`  
    `for s in symbols:`  
        `df1[date, s] = df1[date, s] / df1[0, s]`

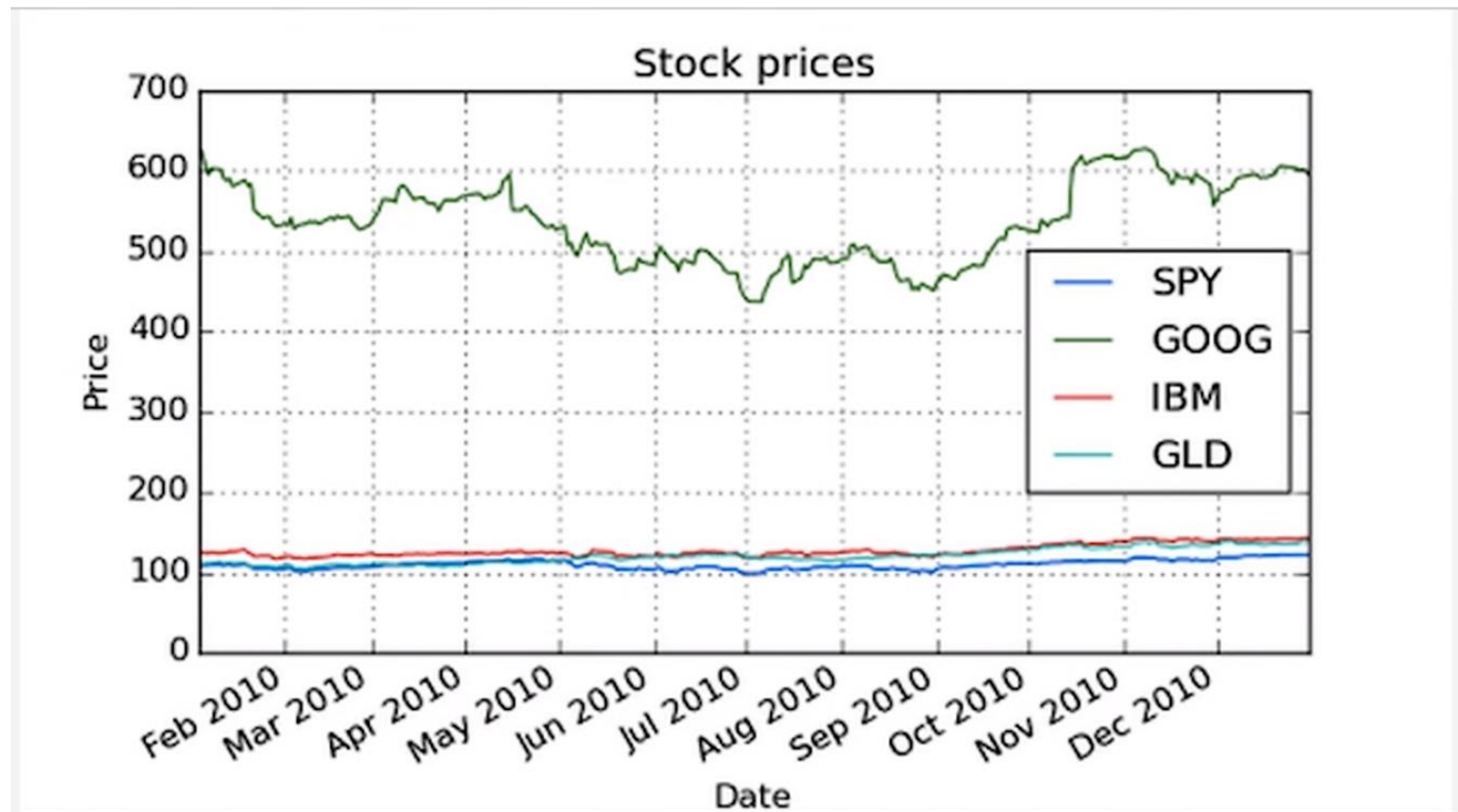
★ `df1 = df1 / df1[0]`

# Plot Function

```
15     if 'SPY' not in symbols: # add SPY for reference, if absent
16         symbols.insert(0, 'SPY')
17
18     for symbol in symbols:
19         df_temp = pd.read_csv(symbol_to_path(symbol), index_col='Date',
20                               parse_dates=True, usecols=['Date', 'Adj Close'], na_values=['nan'])
21         df_temp = df_temp.rename(columns={'Adj Close': symbol})
22         df = df.join(df_temp)
23         if symbol == 'SPY': # drop dates SPY did not trade
24             df = df.dropna(subset=["SPY"])
25
26     return df
27
28 def plot_data(df, title="Stock prices"):
29     '''Plot stock prices'''
30     ax = df.plot(title=title, fontsize=2)
31     ax.set_xlabel("Date")
32     ax.set_ylabel("Price")
33     plt.show() #must be called to show plots in some environments
34
35
36
```



# Plotting multiple stocks





# Quiz: Slice and plot two stocks

```
1 """Slice and plot"""
2
3 import os
4 import pandas as pd
5 import matplotlib.pyplot as plt
6
7
8 def plot_selected(df, columns, start_index, end_index):
9     # TODO: Your code here
10
11     # TODO: Your code here
12     # Note: DO NOT modify anything else!
13
14
15 def symbol_to_path(symbol, base_dir="data"):
16     """Return CSV file path given ticker symbol."""
17     return os.path.join(base_dir, "{}.csv".format(str(symbol)))
18
19
20 def get_data(symbols, dates):
21     """Read stock data (adjusted close) for given symbols from CSV files."""
22     df = pd.DataFrame(index=dates)
23     if 'SPY' not in symbols: # add SPY for reference, if absent
24         symbols.insert(0, 'SPY')
```

**# Slice and plot**

**plot\_selected(df, ['SPY', 'IBM'], '2010-03-01', '2010-04-01')**

```
25 df_temp = df_temp.rename(columns={'Adj Close': symbol})
26 df = df.append(df_temp)
```

# Normalization

```
22         df_temp = df_temp.rename(columns={'Adj Close': symbol})
23         df = df.join(df_temp)
24         if symbol == 'SPY': # drop dates SPY did not trade
25             df = df.dropna(subset=["SPY"])
26
27     return df
28
29
30 def normalize_data(df):
31     """Normalize stock prices using the first row of the dataframe."""
32     return df / df.ix[0,:]
33
34
35 def plot_data(df, title="Stock prices"):
36     """Plot stock prices with a custom title and meaningful axis labels."""
37     ax = df.plot(title=title, fontsize=12)
38     ax.set_xlabel("Date")
39     ax.set_ylabel("Price")
40     plt.show()
41
42
43 def test_normalize():
```

# HAPPY CODING !!

---