

## Computational Finance and FinTech

### File I/O

## Contents

<b>5 Input/Output (IO)</b>	<b>1</b>
5.1 Basic IO with Python	1
5.2 CSV files with <b>pandas</b>	3
5.3 Excel-files with <b>pandas</b>	4
5.4 Saving plots	4
5.5 Direct download from the internet: <b>pandas-datareader</b>	7

## 5 Input/Output (IO)

### Input / Output

- Further reading: **Py4Fi, Chapter 9**
- This chapter is about loading data from and storing data to a persistent data source (e.g. hard drive).
- Often financial data comes as **.csv**-files or **xlsx-** / **xls**-files.
- Before working with the data it must be loaded into Python first.
- Data can be loaded, written and stored using Python's built-in functions or **csv** or **pandas**.
- The book contains additional information on reading from and writing to SQL databases.

### Input / Output

- The usual initialisation:

```
[1]: import numpy as np
import pandas as pd
from pylab import mpl, plt
plt.style.use('seaborn')
%matplotlib inline
```

### 5.1 Basic IO with Python

#### CSV files

- A popular data format for working with spreadsheet-like data is **comma-separated values**, abbreviated as **.csv**.
- In a csv file, data entries are stored in rows, with column separated by comma.
- The top row may contain the column names.

#### CSV files

- Create a dummy data set to demonstrate how to work with csv files:

```
[2]: rows = 10000
a = np.random.standard_normal((rows, 5)).round(4)
a[:2]
```

```
[2]: array([[ 1.4625, -0.0305, -1.4143,  0.3824,  0.1528],
          [-0.4249, -0.244 ,  0.5228,  1.1071, -0.3109]])
```

```
[3]: t = pd.date_range(start='2019/1/1', periods=rows, freq='H') # creates a
      ↪ DateTimeIndex
```

```
[4]: t[:4]
```

```
[4]: DatetimeIndex(['2019-01-01 00:00:00', '2019-01-01 01:00:00',
                    '2019-01-01 02:00:00', '2019-01-01 03:00:00'],
                    dtype='datetime64[ns]', freq='H')
```

Note: 'H' is used for an hourly frequency, see [here](#) for more information about the `freq` settings.

### Creating a .csv-file in Python

- `open()` with the `w` flag creates an empty file.

```
[5]: import os
      path = os.getcwd() + '/' # the current working directory
      path
```

```
[5]: '/Users/nat/Documents/GitHub/compfin_dev/'
```

```
[6]: csv_file = open(path + 'RandomNumbers.csv', 'w') # opens a file for writing
      header = 'date,no1,no2,no3,no4,no5\n'
      csv_file.write(header) # defines the header row and writes it as the first line
```

```
[6]: 25
```

### Creating a .csv-file in Python

- The data is written to the .csv-file using Python's built-in function `write()`:

```
[7]: for t_, (no1, no2, no3, no4, no5) in zip(t, a): # combines the data row-wise ...
      s = '{}{},{}{},{}{},{}{}\n'.format(t_, no1, no2, no3, no4, no5) # ... into string
      ↪ objects ...
      csv_file.write(s) # ... and writes it to the file line-by-line
```

- The changes are saved by calling `close()`.

```
[8]: csv_file.close()
```

```
[9]: ls -l RandomN* # Here it is!
```

```
-rw-r--r--@ 1 nat  staff  569343 23 Mar 20:48 RandomNumbers.csv
-rw-r--r--  1 nat  staff  445582 25 May 2020 RandomNumbers.xlsx
-rw-r--r--  1 nat  staff  579310 25 May 2020 RandomNumbers_2.csv
-rw-r--r--  1 nat  staff  602833 25 May 2020 RandomNumbers_3.csv
```

### Reading a .csv-file in Python

```
[10]: csv_file = open(path + 'RandomNumbers.csv', 'r') # open file for reading ('r')
      RandomNumbers = csv_file.readlines() # read the file contents in one step
      RandomNumbers[:5]
```

```
[10]: ['date,no1,no2,no3,no4,no5\n',
      '2019-01-01 00:00:00,1.4625,-0.0305,-1.4143,0.3824,0.1528\n',
      '2019-01-01 01:00:00,-0.4249,-0.244,0.5228,1.1071,-0.3109\n',
      '2019-01-01 02:00:00,0.7776,0.1785,-0.0359,-0.9896,0.7601\n',
      '2019-01-01 03:00:00,0.3889,0.1908,-0.2494,-1.2765,0.2922\n']
```

```
[11]: csv_file.close() # close the file
```

### csv reader

- CSV files are so common that a dedicated package for reading (and writing) CSV files exists.

```
[12]: import csv

with open(path + 'RandomNumbers.csv', 'r') as f:
    numbers = csv.reader(f) # returns every line as a list object
    lines = [line for line in numbers]

lines[:5]
```

```
[12]: [['date', 'no1', 'no2', 'no3', 'no4', 'no5'],
      ['2019-01-01 00:00:00', '1.4625', '-0.0305', '-1.4143', '0.3824', '0.1528'],
      ['2019-01-01 01:00:00', '-0.4249', '-0.244', '0.5228', '1.1071', '-0.3109'],
      ['2019-01-01 02:00:00', '0.7776', '0.1785', '-0.0359', '-0.9896', '0.7601'],
      ['2019-01-01 03:00:00', '0.3889', '0.1908', '-0.2494', '-1.2765', '0.2922']]
```

```
[13]: type(lines)
```

```
[13]: list
```

**Note:** Aside of a `list`, the `csv` module can also return an `OrderedDict`, by using `DictReader()`. To learn more about the object `OrderedDict`, see [here](#).

### csv writer

- CSV files are easily written to a file using the `csv` package.
- Write the `.csv`-file, using the functions `writer()` and `writerow()`.

```
[14]: with open(path + 'RandomNumbers_2.csv', 'w') as f:
    numbers = csv.writer(f, delimiter = ',')
    for line in lines:
        numbers.writerow(line)
```

## 5.2 CSV files with pandas

- Some of the data formats that `pandas` can read/write:
  - CSV (Comma-separated values)
  - SQL (Structured query language)
  - XLS / XLSX (Microsoft Excel files)
  - JSON (Javascript object notation)
  - HTML (Hypertext markup language)

### CSV files with pandas

- Reading a CSV file with `pandas` loads the data into a `DataFrame` object.
- The first line of the file is assumed to be the header.

```
[15]: df = pd.read_csv('RandomNumbers.csv')
```

```
[16]: df.head()
```

```
[16]:
```

	date	no1	no2	no3	no4	no5
0	2019-01-01 00:00:00	1.4625	-0.0305	-1.4143	0.3824	0.1528
1	2019-01-01 01:00:00	-0.4249	-0.2440	0.5228	1.1071	-0.3109
2	2019-01-01 02:00:00	0.7776	0.1785	-0.0359	-0.9896	0.7601
3	2019-01-01 03:00:00	0.3889	0.1908	-0.2494	-1.2765	0.2922
4	2019-01-01 04:00:00	-0.4529	1.1905	-0.9269	-1.0902	-0.0477

### CSV files with pandas

- `pd.read_csv()` adds an index by default.
- To use a different index, add `set_index(columnName)`.

```
[17]: df.set_index('date', inplace=True)
df.head()
```

```
[17]:
```

	no1	no2	no3	no4	no5
date					
2019-01-01 00:00:00	1.4625	-0.0305	-1.4143	0.3824	0.1528
2019-01-01 01:00:00	-0.4249	-0.2440	0.5228	1.1071	-0.3109
2019-01-01 02:00:00	0.7776	0.1785	-0.0359	-0.9896	0.7601
2019-01-01 03:00:00	0.3889	0.1908	-0.2494	-1.2765	0.2922
2019-01-01 04:00:00	-0.4529	1.1905	-0.9269	-1.0902	-0.0477

### Writing a .csv-file using pandas

```
[18]: df['no1'] = df['no1'] + 1000
df.to_csv('RandomNumbers_3.csv') # couldn't be easier...
```

## 5.3 Excel-files with pandas

- Microsoft Excel files can be read and created easily using pandas.
- This is quite similar to reading / writing CSV files.

```
[19]: df.to_excel('RandomNumbers.xlsx')
```

```
[20]: ls -l RandomNumbers*
```

```
-rw-r--r--@ 1 nat  staff  569343 23 Mar 20:48 RandomNumbers.csv
-rw-r--r--  1 nat  staff  445096 23 Mar 20:48 RandomNumbers.xlsx
-rw-r--r--  1 nat  staff  579344 23 Mar 20:48 RandomNumbers_2.csv
-rw-r--r--  1 nat  staff  602605 23 Mar 20:48 RandomNumbers_3.csv
```

### Reading a .xlsx-file using pandas

```
[21]: df = pd.read_excel('RandomNumbers.xlsx', index_col=0)
```

```
[22]: df.head()
```

```
[22]:
```

	no1	no2	no3	no4	no5
date					

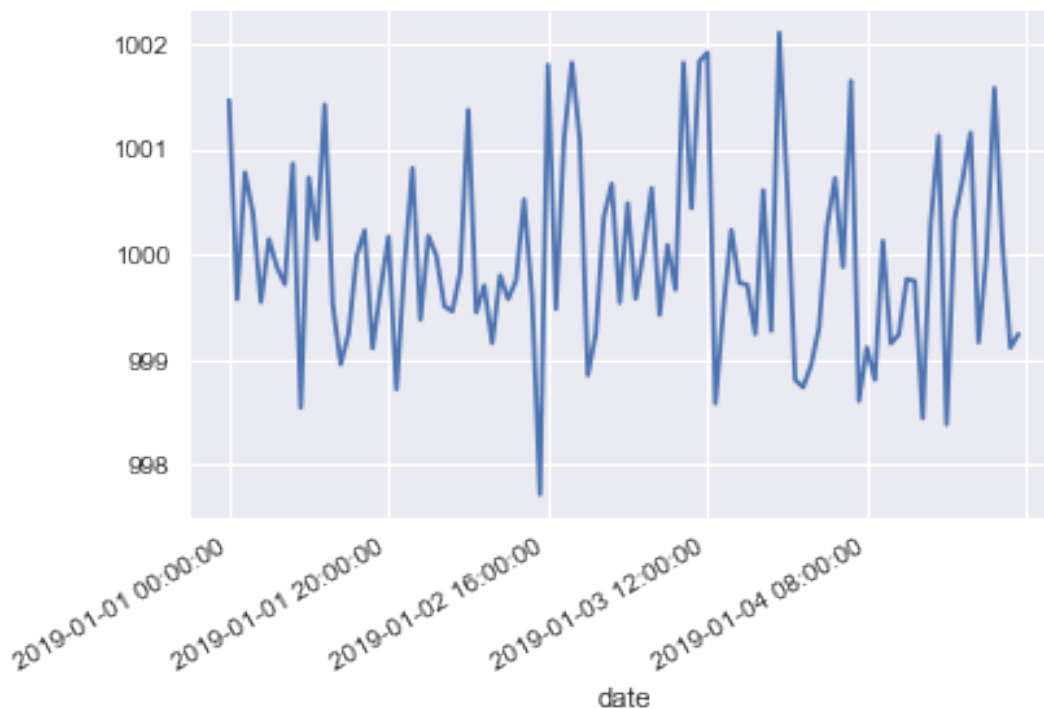
2019-01-01 00:00:00	1001.4625	-0.0305	-1.4143	0.3824	0.1528
2019-01-01 01:00:00	999.5751	-0.2440	0.5228	1.1071	-0.3109
2019-01-01 02:00:00	1000.7776	0.1785	-0.0359	-0.9896	0.7601
2019-01-01 03:00:00	1000.3889	0.1908	-0.2494	-1.2765	0.2922
2019-01-01 04:00:00	999.5471	1.1905	-0.9269	-1.0902	-0.0477

## 5.4 Saving plots

- This section demonstrates how to make plots persistent for further processing.
- The file format is typically detected by the file name extensions.

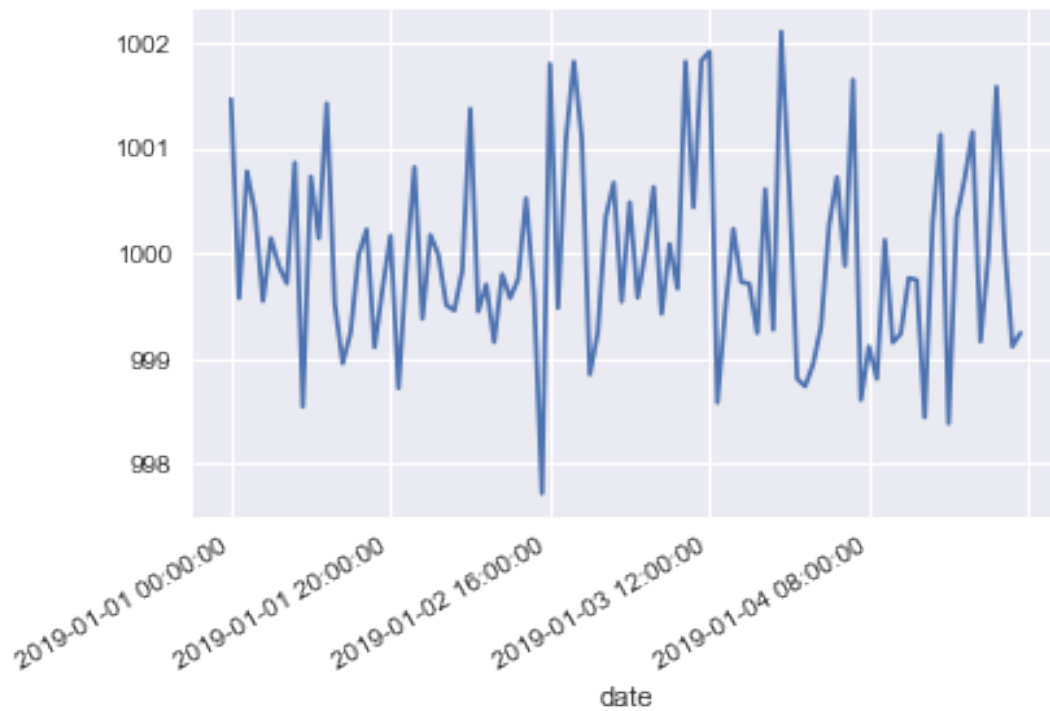
```
[23]: from matplotlib import pyplot as plt
import matplotlib.dates as mdates
```

```
[24]: fig, ax = plt.subplots()
df['no1'][:100].plot()
fig.autofmt_xdate()
plt.savefig('LineChart.png')
```



### Generating a pdf plot

```
[25]: fig, ax = plt.subplots()
df['no1'][:100].plot()
fig.autofmt_xdate()
plt.savefig('LineChart.pdf')
```



### Generating plots

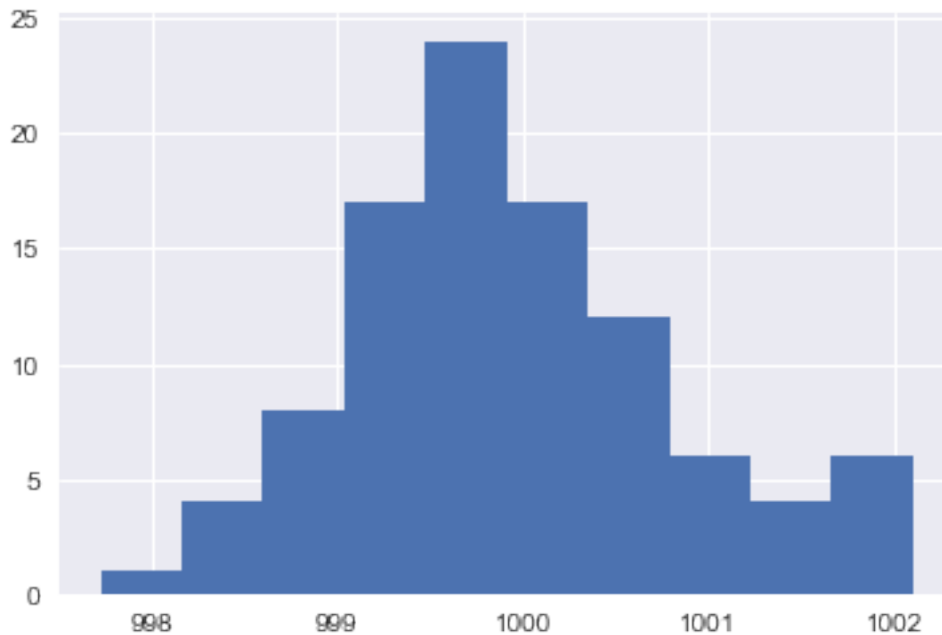
- There are various options for fine-tuning the appearance and quality of the plot.
- Let's check that both plots have been written to the hard disk:

```
[26]: ls -l LineChart*
```

```
-rw-r--r--  1 nat  staff   9851 23 Mar 20:48 LineChart.pdf
-rw-r--r--@ 1 nat  staff 30985 23 Mar 20:48 LineChart.png
```

### Another pdf plot

```
[27]: df['no1'][1:100].hist()
plt.savefig('Hist.pdf')
```



## 5.5 Direct download from the internet: pandas-datareader

- Many data providers offer direct downloads through an API (application programming interface), an programming interface that controls data access and the data format.
- In other words, instead of manually downloading and subsequently importing a file, the data download can directly be incorporated into the program.
- The package `pandas-datareader` provides a convenient unified way of extracting data from various Internet sources into a pandas `DataFrame`.

### Installing pandas-datareader

- Installation via anaconda: go to environments in the left menubar, enter `datareader` in **search packages** and install.
- Alternatively, use a terminal programme (command line) and type:
  - `pip install pandas-datareader`
- This is a frequently changing package... I had to use:
  - `pip install git+https://github.com/pydata/pandas-datareader.git`
- For more information, see [here](#) and [here](#).

### Using pandas-datareader

Using TIINGO requires that you create a free account and obtain an API key.

```
[28]: import os
import pandas_datareader as pdr
```

```
/Users/natalie/anaconda3/lib/python3.7/site-
packages/pandas_datareader/compat/__init__.py:9: FutureWarning:
pandas.util.testing is deprecated. Use the functions in the public API at
pandas.testing instead.
from pandas.util.testing import assert_frame_equal
```

```
[29]: df = pdr.get_data_tingo('GOOG', api_key=os.getenv('TIINGO_API_KEY'))
df.head()
```

```
[29]:
```

			close	high	low	open	volume	\
symbol	date							
GOOG	2016-03-24 00:00:00+00:00		735.30	737.747	731.00	732.01	1594891	
	2016-03-28 00:00:00+00:00		733.53	738.990	732.50	736.79	1301327	
	2016-03-29 00:00:00+00:00		744.77	747.250	728.76	734.59	1903758	
	2016-03-30 00:00:00+00:00		750.53	757.880	748.74	750.10	1782427	
	2016-03-31 00:00:00+00:00		744.95	750.850	740.94	749.25	1718798	

			adjClose	adjHigh	adjLow	adjOpen	\
symbol	date						
GOOG	2016-03-24 00:00:00+00:00		735.30	737.747	731.00	732.01	
	2016-03-28 00:00:00+00:00		733.53	738.990	732.50	736.79	
	2016-03-29 00:00:00+00:00		744.77	747.250	728.76	734.59	
	2016-03-30 00:00:00+00:00		750.53	757.880	748.74	750.10	
	2016-03-31 00:00:00+00:00		744.95	750.850	740.94	749.25	

			adjVolume	divCash	splitFactor
symbol	date				
GOOG	2016-03-24 00:00:00+00:00		1594891	0.0	1.0
	2016-03-28 00:00:00+00:00		1301327	0.0	1.0
	2016-03-29 00:00:00+00:00		1903758	0.0	1.0
	2016-03-30 00:00:00+00:00		1782427	0.0	1.0
	2016-03-31 00:00:00+00:00		1718798	0.0	1.0