

# 파이썬 머신러닝 판다스 데이터분석

## Lecture (1)



Dr. Heesuk Kim

## 강의 교재 소개



출판 : 정보문화사  
저자 : 오승환





## Part 0. 개발환경 준비

## Part 1. 판다스 입문

## Part 2. 데이터 입출력

## Part 3. 데이터 살펴보기

## Part 4. 시각화 도구

## Part 5. 데이터 사전처리

## Part 6. 데이터프레임의 다양한 응용

## Part 7. 머신러닝 데이터 분석





## Part 0. 개발환경 준비

---

1. 아나콘다(Anaconda) 배포판 설치
2. 스파이더(Spyder) 사용법



### 1. 아나콘다(Anaconda) 배포판 설치

- 아나콘다 배포판이란?

- 1) 판다스, 넘파이, 맷플롯립 등 데이터 분석 라이브러리, Spyder 등 개발 도구(IDE)를 통합 지원.
- 2) 버전 관리와 패키지 업데이트가 편리.
- 3) 윈도우, 맥OS, 리눅스 모두 지원.

구분	Anaconda	ActivePython	WinPython
개발자	Anaconda (미국)	ActiveState (캐나다)	WinPython 개발팀
비용	유료/무료	유료/무료	무료(오픈소스)
출시연도	2012년	2006년	2014년
운영체제	윈도우/맥/리눅스	윈도우/맥/리눅스/기타	윈도우
특징	Conda 패키지 관리 그래픽 환경(GUI)	Win32 API 지원 IDLE	WPPM 패키지 관리

[표 0-2] 파이썬 배포판의 종류





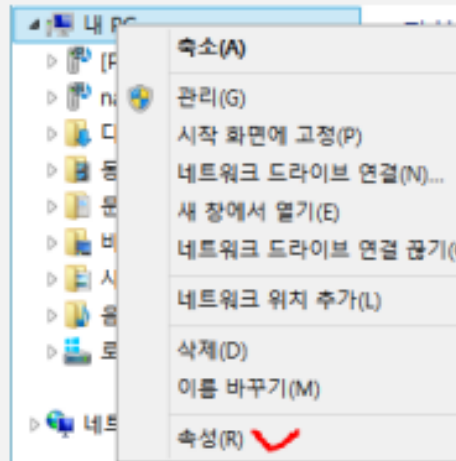
- Checklist before installing Anaconda

- 내 컴퓨터 운영체제 확인: Mac / Windows / Linux

- 운영체제 bit 확인: Mac/Linux → 64bit,

Windows\* → 64bit or 32bit

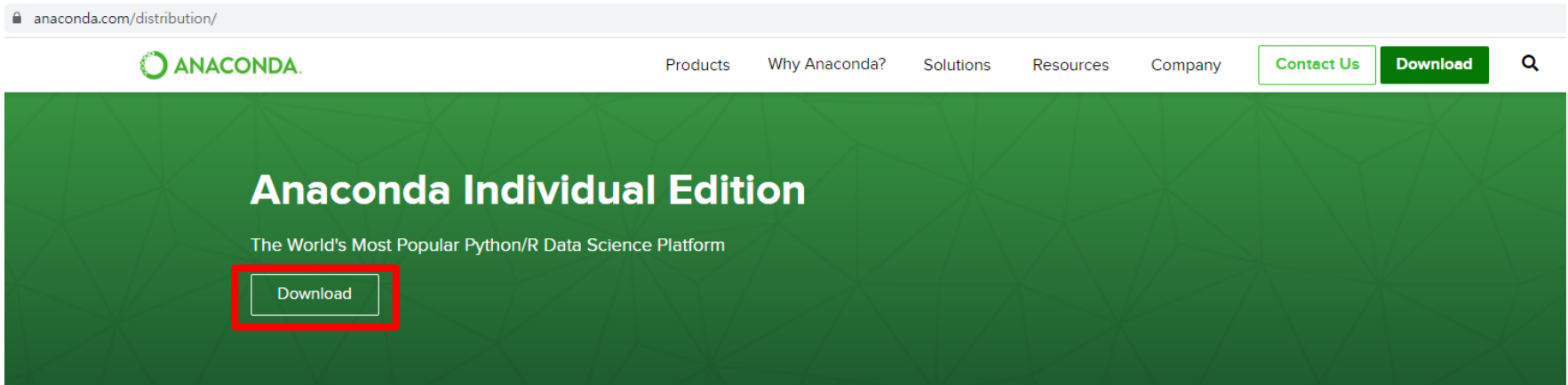
※ 윈도우 탐색기 → 내PC 우클릭 → 속성 → 시스템 종류



# Part 0. 개발환경 준비

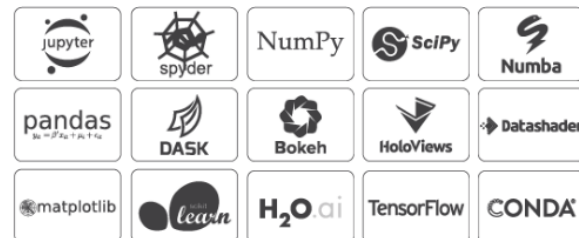


<https://www.anaconda.com/distribution/>



The open-source [Anaconda Individual Edition](#) (formerly Anaconda Distribution) is the easiest way to perform Python/R data science and machine learning on Linux, Windows, and Mac OS X. With over 19 million users worldwide, it is the industry standard for developing, testing, and training on a single machine, enabling *individual data scientists* to:

- Quickly download 7,500+ Python/R data science packages
- Manage libraries, dependencies, and environments with [Conda](#)
- Develop and train machine learning and deep learning models with [scikit-learn](#), [TensorFlow](#), and [Theano](#)
- Analyze data with scalability and performance with [Dask](#), [NumPy](#), [pandas](#), and [Numba](#)
- Visualize results with [Matplotlib](#), [Bokeh](#), [Datashader](#), and [HoloViews](#)



Windows



macOS



Linux





<https://www.anaconda.com/distribution/>



Windows



macOS



Linux

### Anaconda 2020.02 for Windows Installer

#### Python 3.7 version

Download

64-Bit Graphical Installer (466 MB)

32-Bit Graphical Installer (423 MB)

#### Python 2.7 version

Download

64-Bit Graphical Installer (413 MB)

32-Bit Graphical Installer (356 MB)





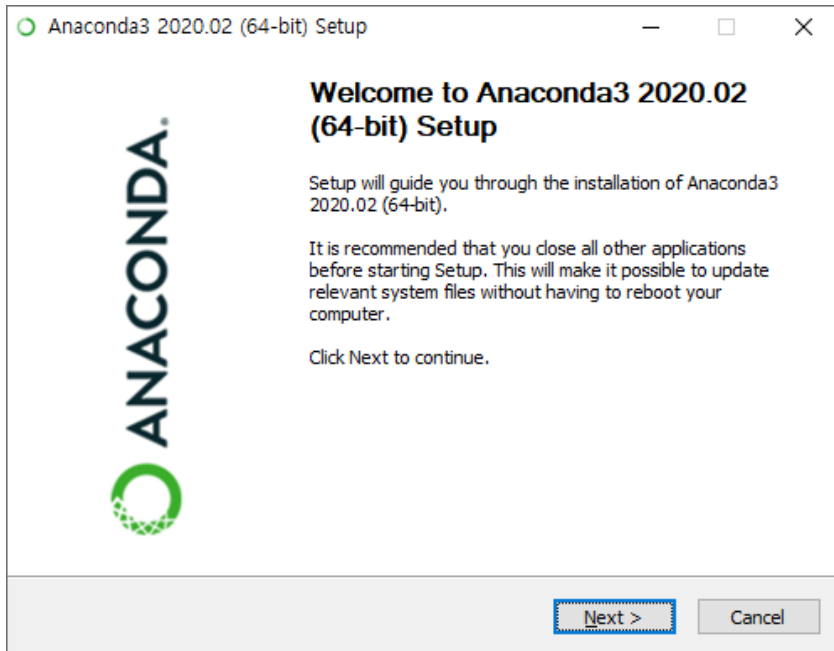
# Part 0. 개발환경 준비



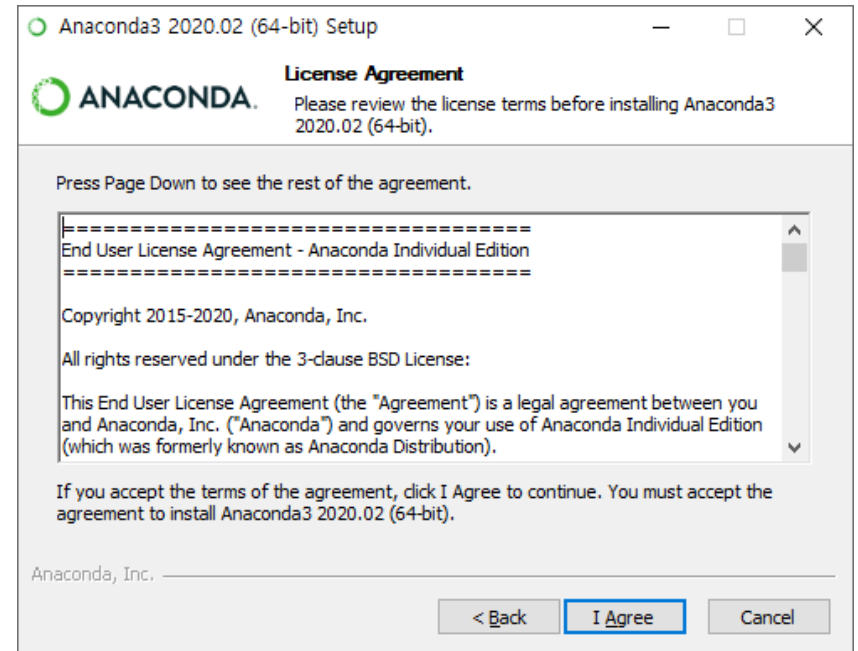
1

<input checked="" type="checkbox"/>  Anaconda3-2020.02-Windows-x86.exe	2020-03-24 오후...	응용 프로그램	477,450KB
---	------------------	---------	-----------

2

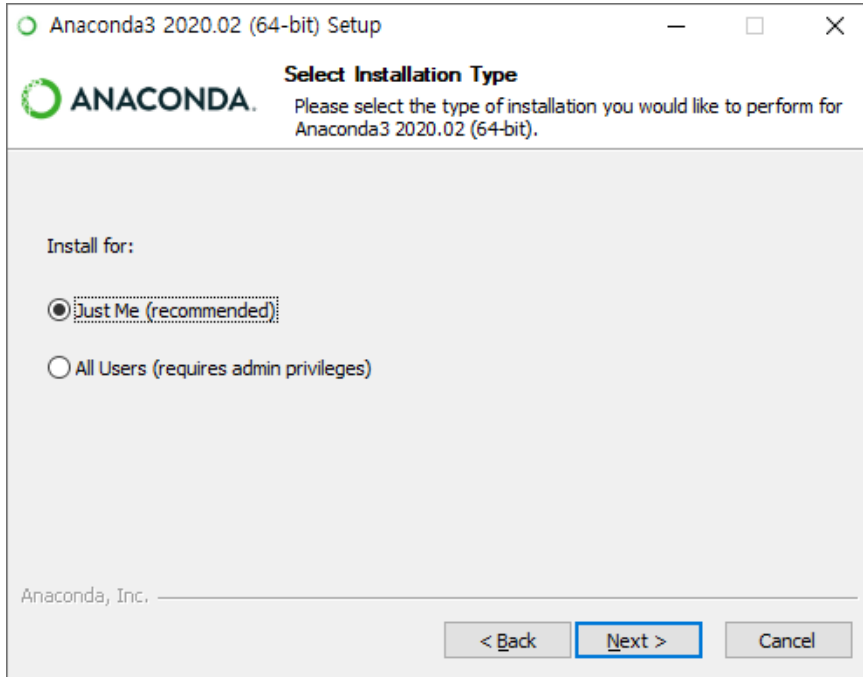


3

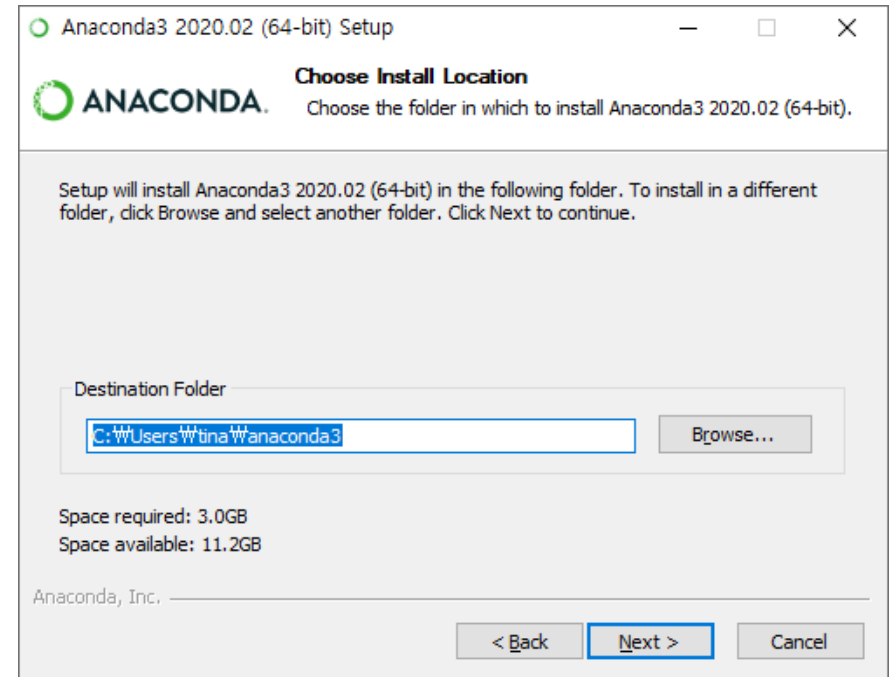




4

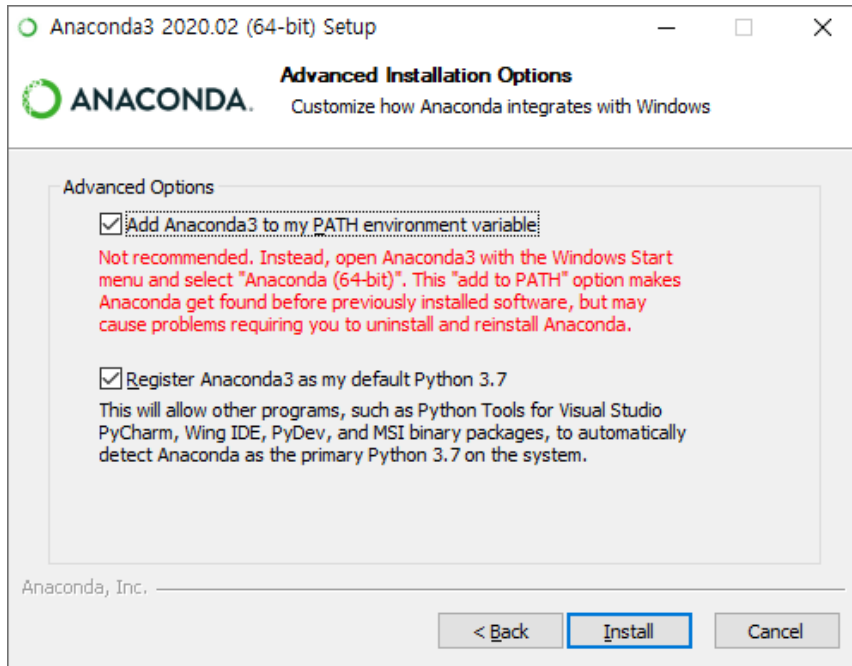


5

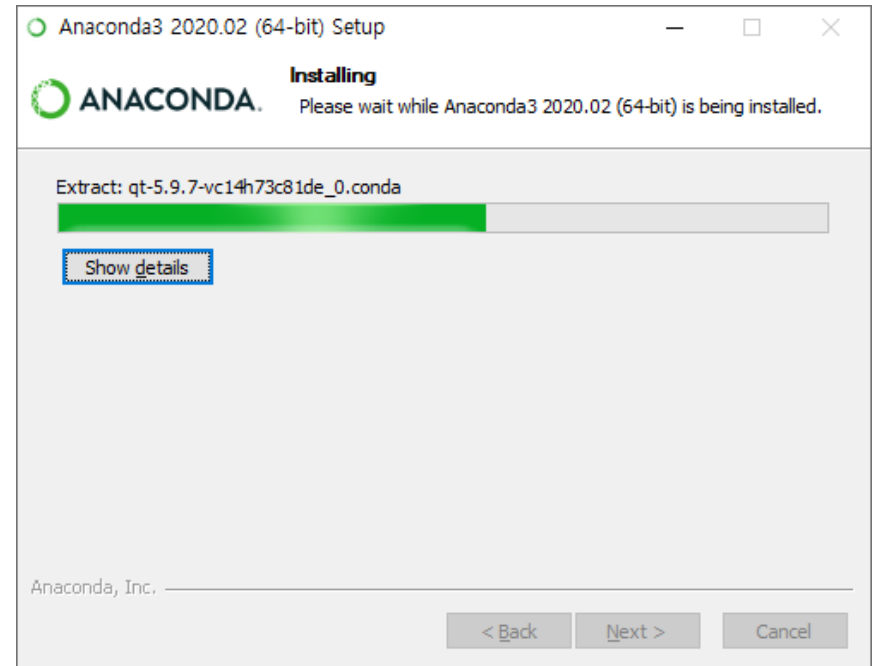




6

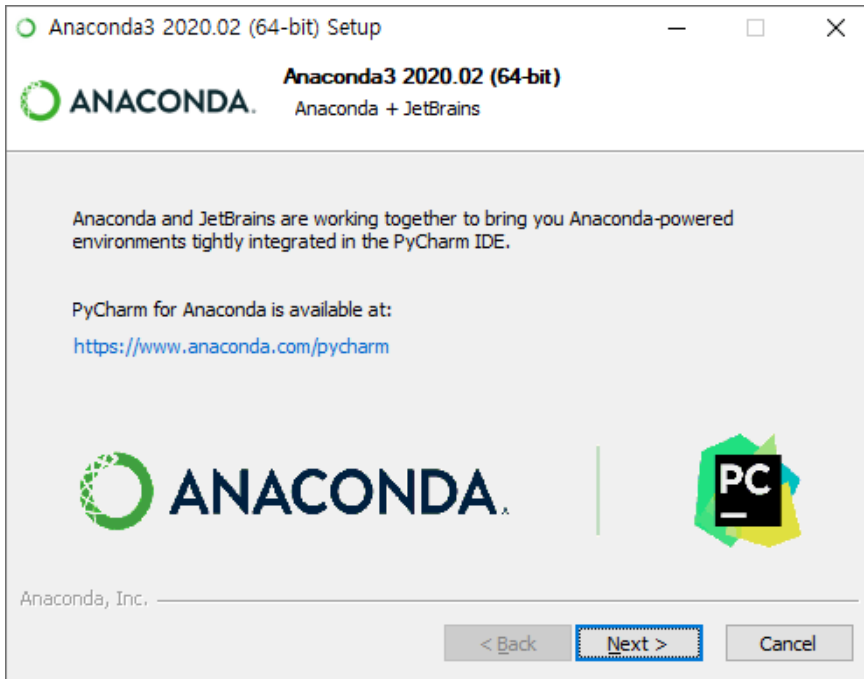


7

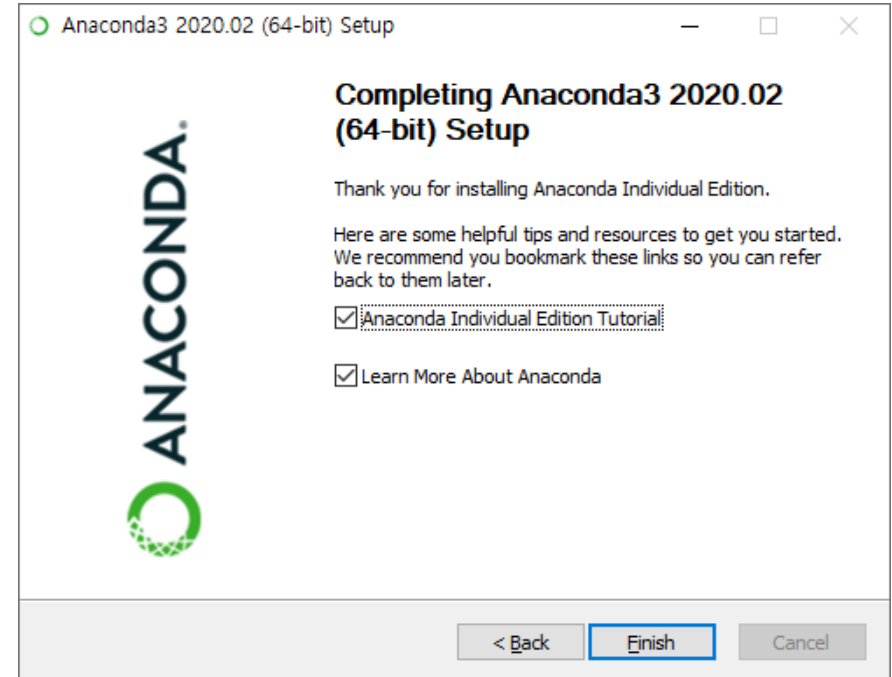




8

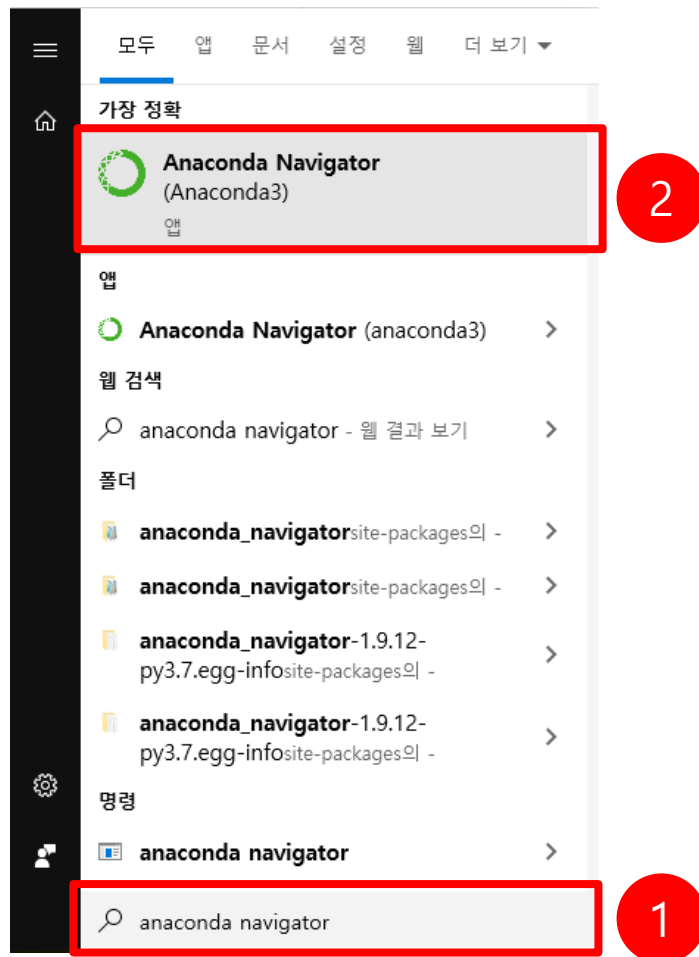


9





- Anaconda install & library 확인



# Part 0. 개발환경 준비



Anaconda Navigator

File Help

ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

1 Home

2 Environments

3 Search Environments

4 base (root)

5 pandas

Name	Description	Version
<input type="checkbox"/> autovizwidget	An auto-visualization library for pandas dataframes	0.13.1
<input type="checkbox"/> blaze	Numpy and pandas interface to big data	0.11.3
<input type="checkbox"/> geopandas	Geographic pandas extensions.	0.6.1
<input type="checkbox"/> pandas	High-performance, easy-to-use data structures and data analysis tools.	1.0.1
<input type="checkbox"/> pandas-datareader	Up to date remote data access for pandas, works for multiple versions of pandas	0.8.1
<input type="checkbox"/> pandas-profiling	Generate profile report for pandas dataframe	1.4.1
<input type="checkbox"/> pandasql	SqlDf for pandas	0.7.3
<input type="checkbox"/> qgrid	Pandas dataframe viewer for jupyter notebook	1.1.1
<input type="checkbox"/> streamz	Manage streaming data, optionally with dask and pandas	0.5.2

9 packages available matching "pandas"

Create Clone Import Remove



# Part 0. 개발환경 준비



Anaconda Navigator

File Help

ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

Home

Environments

Learning

Community

Documentation

Developer Blog

Twitter YouTube GitHub

Create Clone Import Remove

Search Environments

base (root)

All Channels Update index... pandas

Name	Description	Version
<input type="checkbox"/> autovizwidget	An auto-visualization library for pandas dataframes	0.13.1
<input type="checkbox"/> blaze	Numpy and pandas interface to big data	0.11.3
<input type="checkbox"/> geopandas	Geographic pandas extensions.	0.6.1
<input checked="" type="checkbox"/> pandas	High-performance, easy-to-use data structures and data analysis tools.	1.0.1
<input type="checkbox"/> pandas-datareader	Up to date remote data access for pandas, works for multiple versions of pandas	0.8.1
<input type="checkbox"/> pandas-profiling	Generate profile report for pandas dataframe	1.4.1
<input type="checkbox"/> pandasql	SqlDf for pandas	0.7.3
<input type="checkbox"/> qgrid	Pandas dataframe viewer for jupyter notebook	1.1.1
<input type="checkbox"/> streamz	Manage streaming data, optionally with dask and pandas	0.5.2

9 packages available matching "pa[제록 없음]ckage selected

Apply Clear

1

2

# Part 0. 개발환경 준비



Anaconda Navigator

File Help

ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

Home

Environments

Learning

Community

Documentation

Developer Blog

Create Clone Import Remove

Search Environments

base (root)

All Channels Update index...

pandas

Name	Description	Version
<input type="checkbox"/> autovizwidget	An auto-visualization library for pandas dataframes	0.13.1
<input type="checkbox"/> blaze	Numpy and pandas interface to big data	0.11.3
<input type="checkbox"/> geopandas	Geographic pandas extensions.	0.6.1
<input checked="" type="checkbox"/> pandas	High-performance, easy-to-use data structures and data analysis tools.	<a href="#">0.24.2</a>
<input type="checkbox"/> pandas-datareader	Up to date remote data access for pandas, works for multiple versions of pandas	0.8.1
<input type="checkbox"/> pandas-profiling	Generate profile report for pandas dataframe	1.4.1
<input type="checkbox"/> pandasql	SqlDf for pandas	0.7.3
<input type="checkbox"/> qgrid	Pandas dataframe viewer for jupyter notebook	1.1.1
<input type="checkbox"/> streamz	Manage streaming data, optionally with dask and pandas	0.5.2

9 packages available matching "pandas"





# Part 0. 개발환경 준비



Anaconda Navigator

File Help

ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

Home

Environments

Learning

Community

Documentation

Developer Blog

Create Clone Import Remove

Search Environments

base (root)

All Channels Update index...

numpy

Name	Description	Version
<input type="checkbox"/> blaze	Numpy and pandas interface to big data	0.11.3
<input type="checkbox"/> bottleneck	Fast numpy array functions specialized for use in orange	0.7.1
<input checked="" type="checkbox"/> bottleneck	Fast numpy array functions written in cython.	1.2.1
<input type="checkbox"/> cupy	Cupy is an implementation of a numpy-compatible multi-dimensional array on cuda.	6.0.0
<input checked="" type="checkbox"/> mkl_fft	Numpy-based implementation of fast Fourier transform using intel (r) math kernel library.	1.0.12
<input checked="" type="checkbox"/> mkl_random	Intel (r) mkl-powered package for sampling from common probability distributions into numpy arrays.	1.0.2
<input type="checkbox"/> msgpack-numpy	Numpy data serialization using msgpack	0.4.4.3
<input checked="" type="checkbox"/> numba	Numpy aware dynamic python compiler using llvm	0.44.1
<input checked="" type="checkbox"/> numexpr	Fast numerical expression evaluator for numpy.	2.6.9
<input checked="" type="checkbox"/> numpy	Array processing for numbers, strings, records, and objects.	1.17.2
<input checked="" type="checkbox"/> numpy-base		1.17.2
<input type="checkbox"/> numpy-devel		1.9.3
<input checked="" type="checkbox"/> numpydoc	Sphinx extension to support docstrings in numpy format	0.9.1

15 packages available matching "numpy"



# Part 0. 개발환경 준비



Anaconda Navigator

File Help

ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

Home

Environments

Learning

Community

Documentation

Developer Blog

Twitter YouTube GitHub

Create Clone Import Remove

Search Environments

base (root)

All Channels Update index... matplotlib X

Name	Description	Version
<input type="checkbox"/> basemap	Plot on map projections using matplotlib	1.2.0
<input type="checkbox"/> descartes	Use geometric objects as matplotlib paths and patches.	1.1.0
<input checked="" type="checkbox"/> matplotlib	Publication quality figures in python	<a href="#">3.1.0</a>
<input type="checkbox"/> matplotlib-base		3.1.3
<input type="checkbox"/> mpl-scatter-density	Matplotlib helpers to make density scatter plots	0.6
<input type="checkbox"/> mpld3	D3 viewer for matplotlib.	0.3

6 packages available matching "matplotlib"



# Part 0. 개발환경 준비



Anaconda Navigator

File Help

ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

Home

Environments

Learning

Community

Documentation

Developer Blog

Twitter YouTube GitHub

Create Clone Import Remove

Search Environments

base (root)

All Channels Update index...

scipy

Name	Description	Version
✓ scikit-image	Image processing routines for scipy.	0.15.0
✓ scipy	Scientific library for python	1.2.1
✓ statsmodels	Statistical computations and models for use with scipy	0.10.0

3 packages available matching "scipy"



# Part 0. 개발환경 준비



Anaconda Navigator

File Help

ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

Home

Environments

Learning

Community

Documentation

Developer Blog

Twitter YouTube GitHub

Create Clone Import Remove

Search Environments

base (root)

All Channels Update index...

scikit

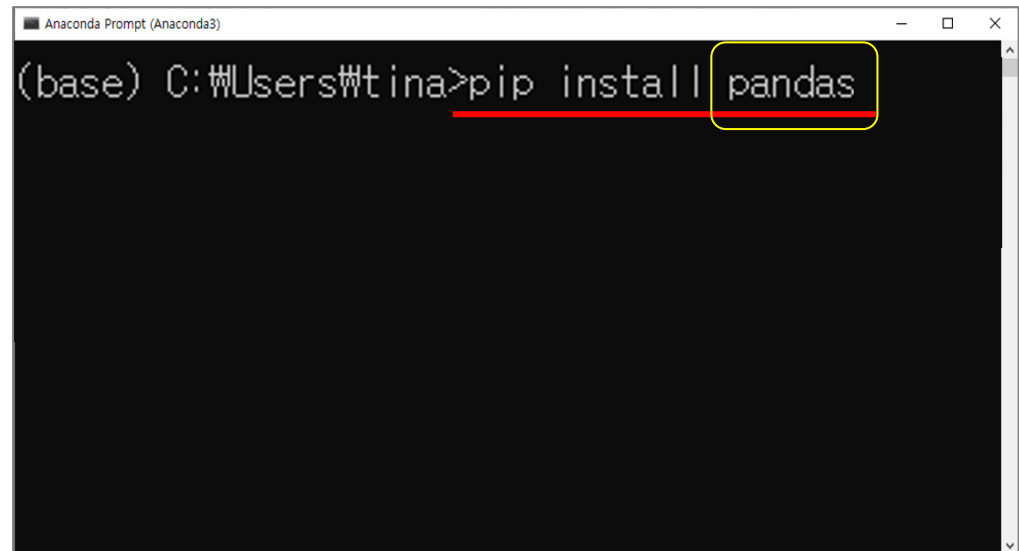
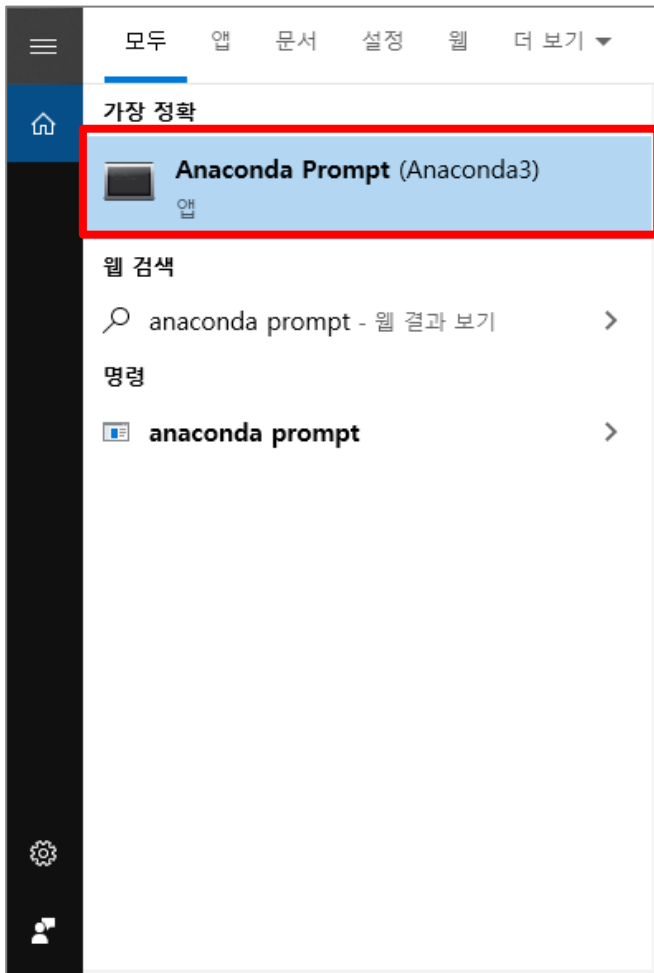
Name	Description	Version
<input type="checkbox"/> dask-searchcv	Tools for doing hyperparameter search with scikit-learn and dask	0.2.0
<input checked="" type="checkbox"/> scikit-image	Image processing routines for scipy.	<a href="#">0.15.0</a>
<input checked="" type="checkbox"/> scikit-learn	A set of python modules for machine learning and data mining	<a href="#">0.21.2</a>
<input type="checkbox"/> scikit-rf	Object oriented microwave engineering.	0.14.9

4 packages available matching "scikit"





참고



```
pip install pandas  
pip install numpy  
pip install matplotlib  
pip install scipy  
pip install -U scikit-learn
```



# Part 0. 개발환경 준비



ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

[Home](#)  
[Environments](#)  
[Learning](#)  
[Community](#)

Applications on base (root) Channels Refresh

CMD.exe Prompt  
0.1.1  
Run a cmd.exe terminal with your current environment from Navigator activated  
[Launch](#)

JupyterLab  
1.2.6  
An extensible environment for interactive and reproducible computing, based on the Jupyter Notebook and Architecture.  
[Launch](#)

Notebook  
6.0.3  
Web-based, interactive computing notebook environment. Edit and run human-readable docs while describing the data analysis.  
[Launch](#)

Powershell Prompt  
0.0.1  
Run a Powershell terminal with your current environment from Navigator activated  
[Launch](#)

PyCharm  
2019.3.4  
Full-featured Python IDE by JetBrains. Supports code completion, linting, debugging, and domain-specific enhancements for web development and data science.  
[Launch](#)

Qt Console  
4.6.0  
PyQt GUI that supports inline figures, proper multiline editing with syntax highlighting, graphical calltips, and more.  
[Launch](#)

Spyder  
4.0.1  
Scientific PYTHON Development Environment. Powerful Python IDE with advanced editing, interactive testing, debugging and introspection Features  
[Launch](#)

Glueviz  
0.15.2  
Multidimensional data visualization across files. Explore relationships within and among related datasets.  
[Install](#)

Orange 3  
3.23.1  
Component based data mining framework. Data visualization and data analysis for novice and expert. Interactive workflows with a large toolbox.  
[Launch](#)

RStudio  
1.1.456  
A set of integrated tools designed to help you be more productive with R. Includes R essentials and notebooks.  
[Launch](#)

Documentation  
Developer Blog

# Part 0. 개발환경 준비



ANACONDA NAVIGATOR

Sign in to Anaconda Cloud

[Home](#)  
[Environments](#)  
[Learning](#)  
[Community](#)

Applications on base (root) Channels Refresh

CMD.exe Prompt  
0.1.1  
Run a cmd.exe terminal with your current environment from Navigator activated  
[Launch](#)

JupyterLab  
1.2.6  
An extensible environment for interactive and reproducible computing, based on the Jupyter Notebook and Architecture.  
[Launch](#)

Notebook  
6.0.3  
Web-based, interactive computing notebook environment. Edit and run human-readable docs while describing the data analysis.  
[Launch](#)

Powershell Prompt  
0.0.1  
Run a Powershell terminal with your current environment from Navigator activated  
[Launch](#)

PyCharm  
2019.3.4  
Full-featured Python IDE by JetBrains. Supports code completion, linting, debugging, and domain-specific enhancements for web development and data science.  
[Launch](#)

Qt Console  
4.6.0  
PyQt GUI that supports inline figures, proper multiline editing with syntax highlighting, graphical calltips, and more.  
[Launch](#)

Spyder  
4.0.1  
Scientific PYTHON Development Environment. Powerful Python IDE with advanced editing, interactive testing, debugging and introspection Features  
[Launch](#)

Glueviz  
0.15.2  
Multidimensional data visualization across files. Explore relationships within and among related datasets.  
[Install](#)

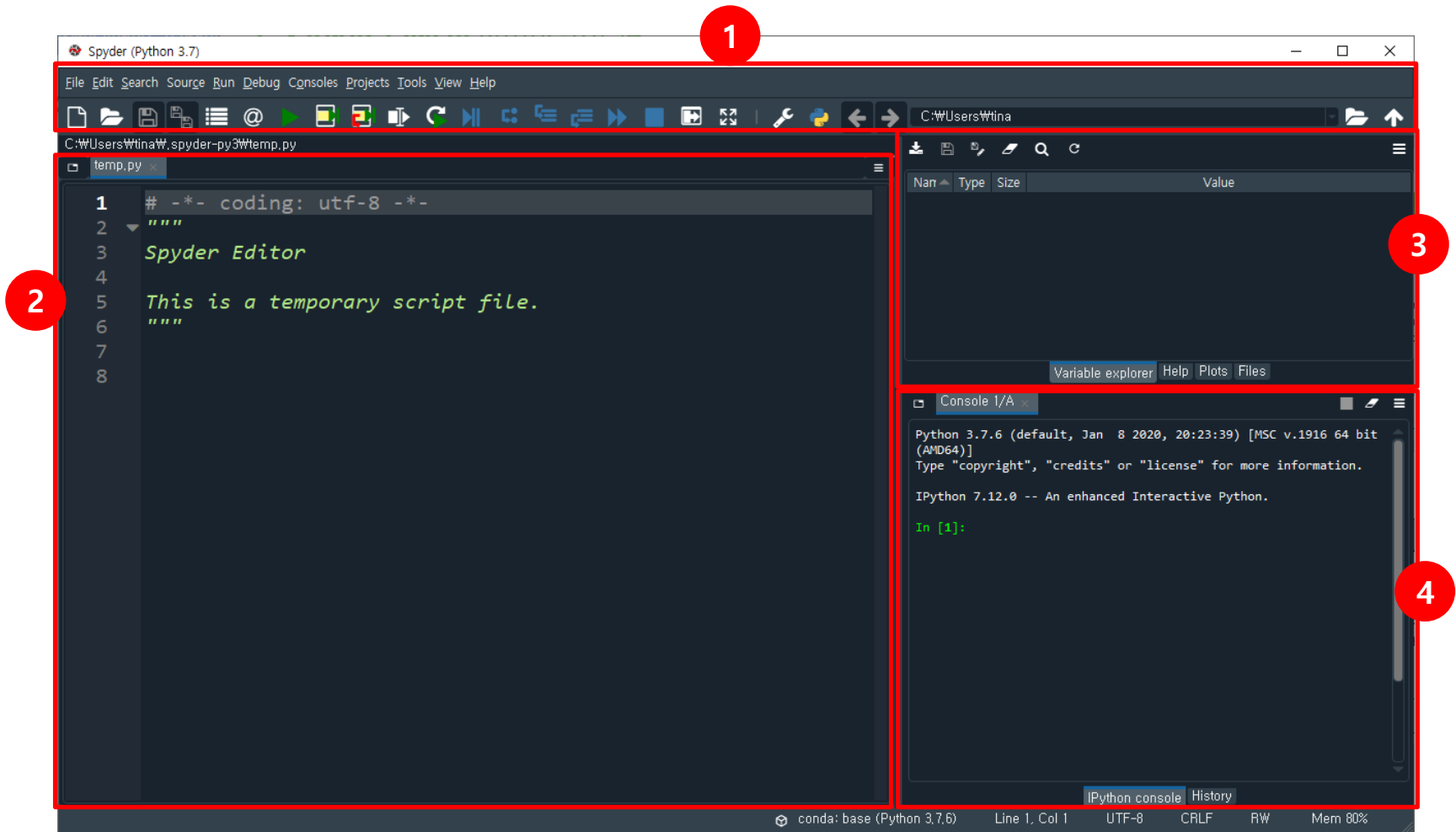
Orange 3  
3.23.1  
Component based data mining framework. Data visualization and data analysis for novice and expert. Interactive workflows with a large toolbox.

RStudio  
1.1.456  
A set of integrated tools designed to help you be more productive with R. Includes R essentials and notebooks.

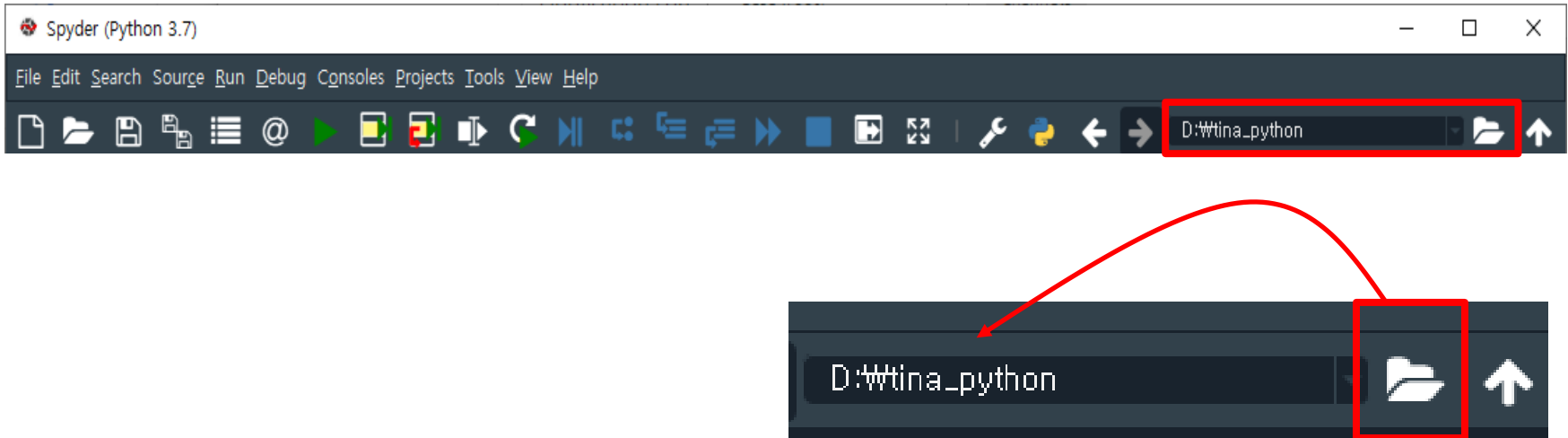
[Documentation](#)  
[Developer Blog](#)



## 2. 스파이더(Spyder) 사용법







### Homework

작업 폴더 밑에 "**본인영어이름\_python**"으로 폴더 만든 후  
Spyder 화면 캡처하여 제출

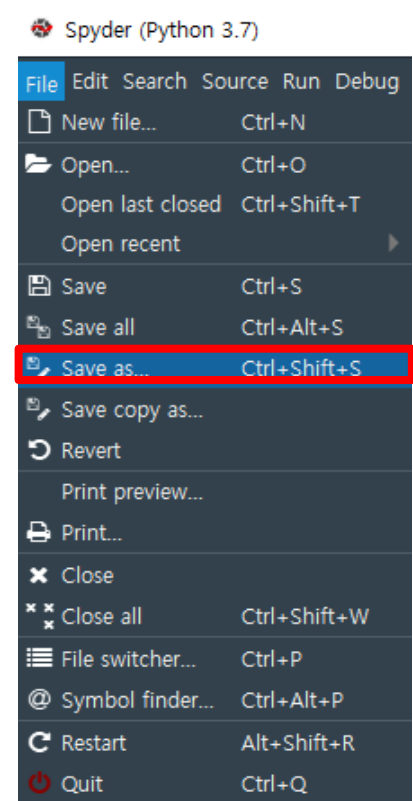
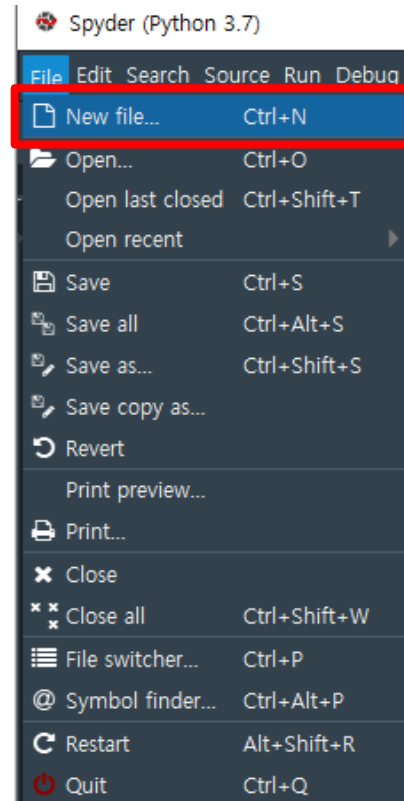
예) 홍길동

*C:\#users#gildong\_python 또는 D:\#users#gildong\_python*





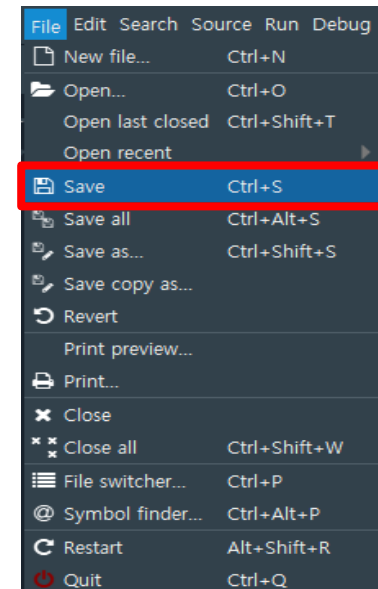
[File] - [New File]을 선택하면 새로운 파일이 만들어진다. 좌측에 있는 에디터 창에 입력한 코드를 실행하려면, 코드 작성 전에 반드시 파일명을 지정하여 저장해야 한다. [File] - [Save As]를 선택하고, 폴더를 지정하여 'sample.py'라는 이름으로 저장한다. 에디터 창에 다음과 같이 코드를 입력하고, 상단 메뉴에서 [File] - [Save]를 선택하거나 저장 버튼(💾)을 클릭하면 파일에 저장된다.





[File] - [New File]을 선택하면 새로운 파일이 만들어진다. 좌측에 있는 에디터 창에 입력한 코드를 실행하려면, 코드 작성 전에 반드시 파일명을 지정하여 저장해야 한다. [File] - [Save As]를 선택하고, 폴더를 지정하여 'sample.py'라는 이름으로 저장한다. 에디터 창에 다음과 같이 코드를 입력하고, 상단 메뉴에서 [File] - [Save]를 선택하거나 저장 버튼(💾)을 클릭하면 파일에 저장된다.

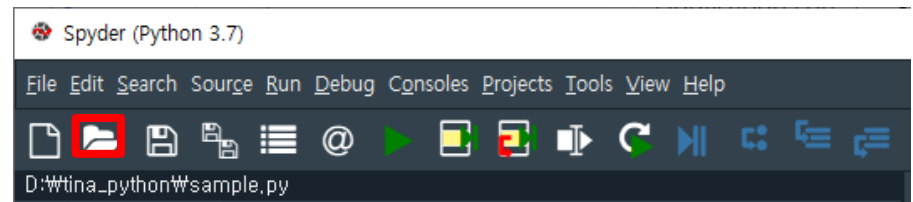
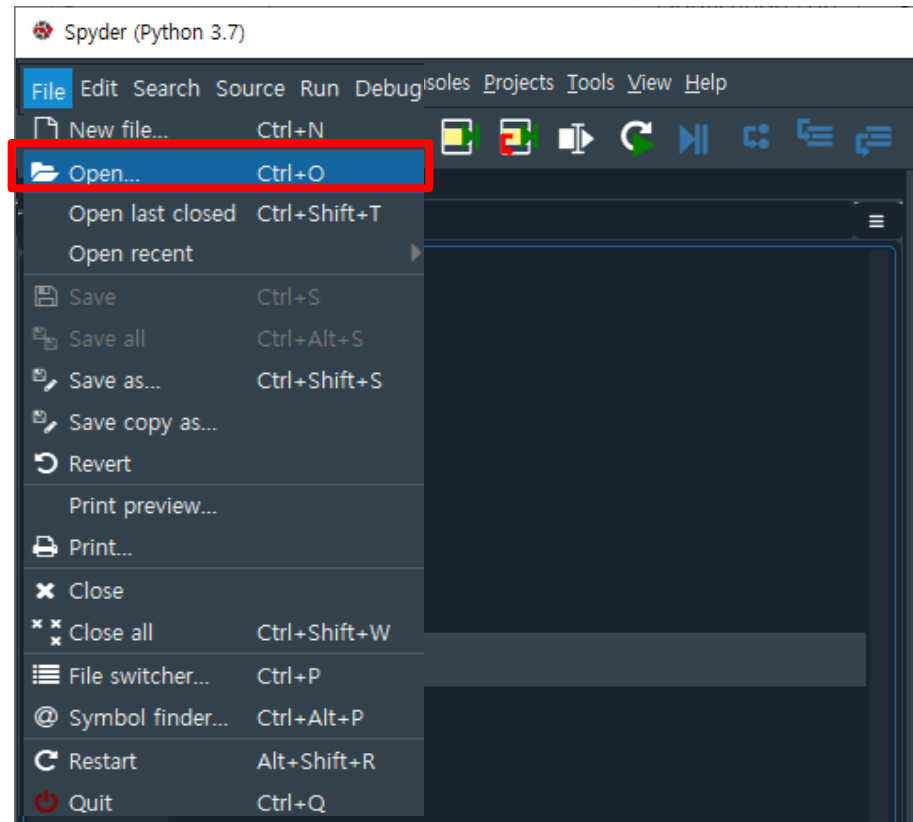
```
1 a = 5
2 print(a)
3
4 a = a+5
5
6 print(a)
7
8
9
```



## Part 0. 개발환경 준비



로컬 PC에 저장되어 있는 파이썬 파일을 열어서 에디터 창에서 편집하려면, 상단 메뉴 [File] - [Open]을 선택하거나 열기 버튼(📁)을 클릭한다. 또는 주소창 옆의 폴더 검색 버튼(🔍)을 눌러서 폴더를 지정하고, 'File Explorer' 탭을 클릭하면 'sample.py'라는 파일명을 찾을 수 있다. 파일명을 더블클릭하면 파일이 열리면서 에디터 창이 활성화된다.





에디터 창에 입력하고 저장한 예제 코드 ('sample.py') 전체를 한번에 실행하려면, 상단 메뉴의 전부 실행 버튼(▶)을 클릭하거나 단축 키 [F5]를 누른다. 우측 하단의 IPython 콘솔에 다음과 같이 실행 결과가 출력된다(변수 a에 저장된 값은 실행 순서에 따라 첫 줄에 5가 출력되고, 두 번째 줄에 다시 5가 더해진 10이 출력된다).

```
Python 3.7.6 (default, Jan 8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)]
Type "copyright", "credits" or "license()" for more information.

IPython 7.12.0 -- An enhanced Interactive Python.

In [1]: runfile('C:/Users/tina/untitled3.py', wdir='C:/Users/tina')
5
10

In [2]: runfile('C:/Users/tina/untitled3.py', wdir='C:/Users/tina')
5
10

In [3]: runfile('D:/python_prog_ex/test_1.py', wdir='D:/python_prog_ex')
5
10

In [4]: runfile('D:/tina_python/sample.py', wdir='D:/tina_python')
5
10

In [5]:
```

IPython console History

conda: base (Python 3.7.6) Line 8, Col 1 UTF-8 CRLF RW Mem 80%





### (IPython) 콘솔 모드

IPython 콘솔은 인터랙티브 실행 환경이다. 따라서 콘솔에 표시된 명령 프롬프트에 파이썬 코드를 직접 입력하면 바로 실행 결과가 출력된다. 다음 그림에서 “In [4]:” 오른쪽에 커서를 놓고 원하는 코드를 입력하면 ‘Out [4]:’ 우측에 실행 결과가 출력된다. 변수

a(현재값 5)에 숫자 5를 더하는 식을 입력한 결과, 숫자 10이 출력된 것을 확인할 수 있다. 다음 줄에는 다음의 명령 입력을 위해 ‘In [5]:’와 같이 명령 프롬프트가 다시 나타난다.

```
Console 4/A x
In [3]: a
Out[3]: 5

In [4]: a + 5
Out[4]: 10

In [5]: |
```







# Any Question?

[edishskim@naver.com](mailto:edishskim@naver.com)

Thank you.

