# 2019-01-17 Lab 1

## Contents

# Exercise 1: Install Anaconda

☐ Go to the Anaconda download page, find the distribution for your operating system, and download it.

☐ Install it once it is downloaded.

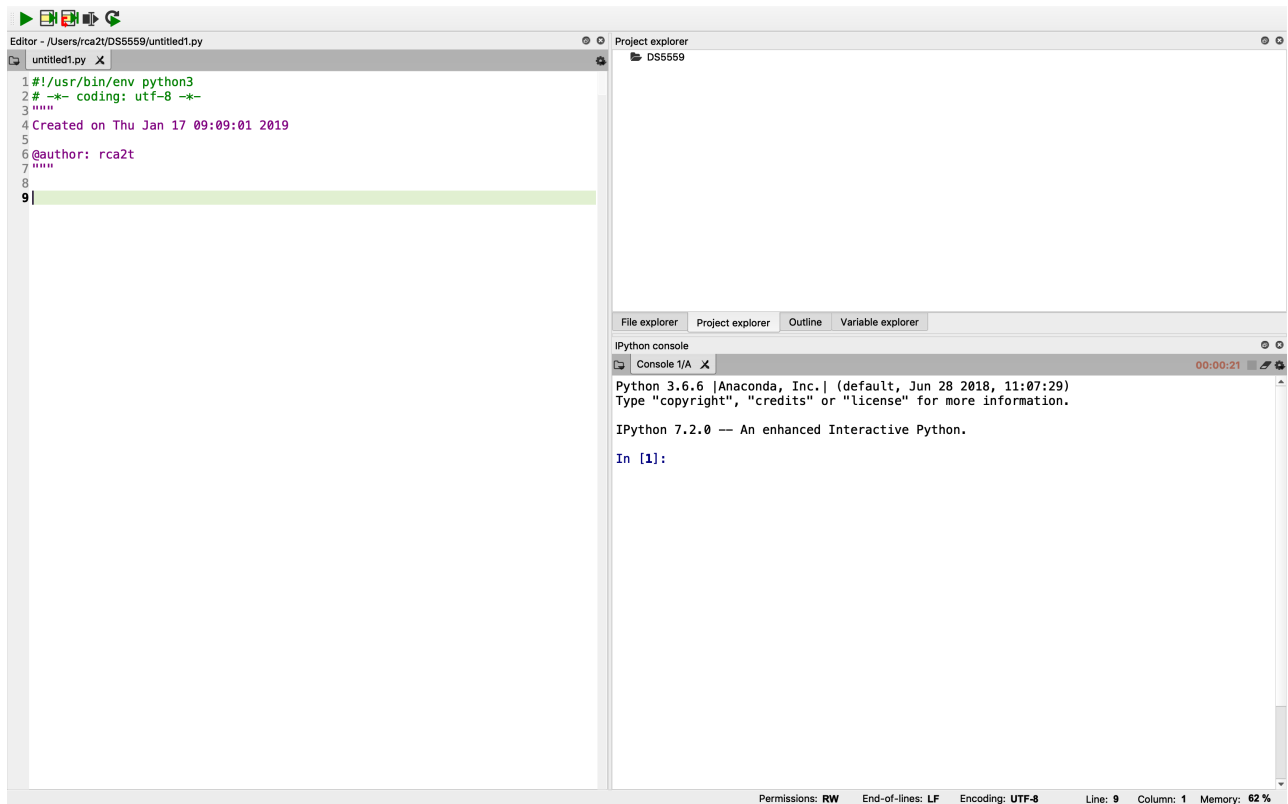☐ Find the application icon in your toolbar or program manager and open it up – you should see something like this:



# Exercise 2: Install Spyder and JupyterLabs

☐ In Anaconda Navigator, locate the cards for Jupyterlab and Spyder and click the Install button for each. You may also install Jupyter Notebook if you'd like.

☐ Once installed, click the Launch button to run it.

☐ Spyder should appear.
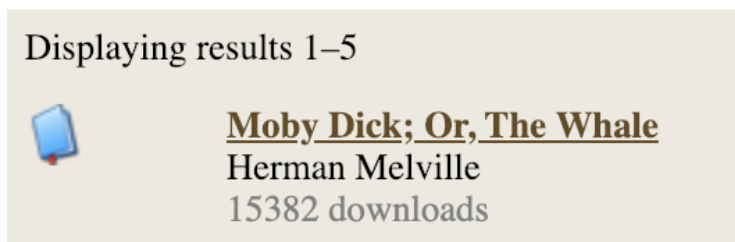
☐ Optionally configure the panels as instructed.

# Exercise 3: Create a New Project in Spyder for this Course

☐ In the menu, select **Projects New Project**.

☐ Select **New Directory** and create a project called **DS5559**.

☐ You should see something like this:



# Exercise 4: Download and Open *Moby Dick* from Project Gutenberg

☐ Go to Project Gutenberg.

☐ Entry "Moby Dick" in the search box on the left.

☐ Click on the first result –



☐ Click on the link for the Plaintext UTF-8 version and save the file to your computer. Save the file to your project directory or move the file there once it is downloaded.

☐ From within Spyder, open the file.

☐ Investigate the file along with the class.

    ☐ Is the file complete?

    ☐ How many lines does the file have?

    ☐ What kind of text encoding does it have?

    ☐ What kind of line endings does it have?

☐ Does the file come with metadata?

☐ What content do we want to keep?

☐ Where does the cruft begin and end?

# Exercise 5: Open File in Python and Convert to Text From 1

☐ Create a new Python file called `moby.py`.

☐ Follow instructor for specific commands.

# Code

```
# -*- coding: utf-8 -*-

import re
import pandas as pd


# Identify the source text (F0)
src_file = '2701-0.txt'

# Import the text as list of lines
lines = open(src_file, 'r', encoding='utf-8').readlines()

# Trim the cruft we identified
lines = lines[340:21964]

# Convert the lines into one big line, preserving line breaks
bigline = ''.join(lines)

# Split the bigline into paragraphs
paras = re.split(r'\n\n+', bigline)

# Break line into paragraphs

# Split by non-character
paras2 = []
for para in paras:
    tokens = re.split(r'\W+', para)
    paras2.append(tokens)


# Split by non-character but keep them
paras3 = []
for para in paras:
    tokens = re.split(r'(\W+)', para)
    paras3.append(tokens)

# Split by non-character using list comprehension
```

```python
paras4 = [re.split(r'\W+', para) for para in paras]

# Try in Pandas

# Import paragraphs into a data frame
df = pd.DataFrame(paras, columns=['line'])
df.index.name = 'line_id'
df.line = df.line.str.strip()

# Tokenize using this one trick
df2 = df.line.str.split(r'\W+', expand=True)\
    .stack()\
    .to_frame()\
    .rename(columns={0: 'token'})
df2.index.names = ['line_id', 'token_id']

# Do a simple normalization
df2['norm'] = df2.token.str.lower()

# Get top N tokens
N = 30
df2['norm'].value_counts().head(N).sort_values().plot(kind='barh')
```

```python
# Visualize dispersion plots of 'ahab' and 'whale'
(df2['norm'] == 'whale').astype('int').plot(figsize=(10, 1))
(df2['norm'] == 'ahab').astype('int').plot(figsize=(10, 1))
```