

Tutorial Pack 9

(90 minutes)

(To be completed during LW9 tutorial)

LEARNING OBJECTIVES

- To revise core Python concepts and practice reading Python statements.
- To develop a Python program that communicates with the Reddit API.

LEARNING OUTCOMES

- By the end of this tutorial students will have:
 - Revised key Python programming concepts
 - Familiarised themselves with object-orientated programming principles
 - Setup a Reddit developer account and obtained access keys.
 - Practiced collecting data using the Reddit API
 - Practiced using external Python packages (Pandas)

RESOURCES AND TOOLS REQUIRED

- Python 3 via google collab or repl.it

IMPORTANT:

This pack is designed for you to go at your own speed. At the end of each section there is a series of practice questions and exercises. You should attempt to answer **all questions**.

Any questions you do not complete today should be completed before your next tutorial.

You may find this tutorial pack and the exercises useful preparation for the coursework.

1. Introduction

In LW8 we covered the basics of Python programming. We will begin with a revision of last week's material and a review of some important Python concepts.

2. Revision of Python concepts (30mins)

Using your understanding of how Python programs work answer the following questions. You should consult your notes from last week's tutorial if you need help.

QNum	Question
1	What will be displayed when the following Python code is executed? print("windy and overcast weather".title())
	Answer
QNum	Question
2	What will be displayed when the following Python code is executed? x = 2 y = 3 print(x**y)
	Answer
QNum	Question
3	What will be displayed when the following Python code is executed? i = 10 if i > 5: print(i*2) else: print(i*3)
	Answer

QNum	Question
4	<p>What will be displayed when the following Python code is executed?</p> <pre>k = [1,2,3,4,5] for i in k: print(i + i)</pre>
	Answer
QNum	Question
5	<p>What will be displayed when the following Python code is executed?</p> <pre>tweet = "Today is a new day!" print(tweet[-4])</pre>
	Answer
QNum	Question
6	<p>What will be displayed when the following Python code is executed?</p> <pre>tweet = "Today is a new day!" print(tweet.split()[3])</pre>
	Answer
QNum	Question
7	<p>Write a python statement to print the number of characters in the following tweet</p> <pre>tweet = "I've just watched the film Gravity...What's up with that? At one point, Scarlett Johansson was in the running to star in this film, as was Naomi Watts"</pre>
	Answer

QNum	Question
8	<p>Explain the purpose of the following program. What will be the output when it is run?</p> <pre>def stopwords(text): found = [] for word in text.split(): if word in ["the", "from", "to", "at", "they", "them"]: if word not in found: found.append(word) return found print(stopwords("Make a donation, help save lives around the world"))</pre>
	Answer
QNum	Question
9	<p>Why does the following program fail to run? What changes should you make to fix it?</p> <pre>class Tweet: def __init__(self, author, text, post_date): self._author = author self._text = text self._post_date = post_date user_tweet = Tweet()</pre>
	Answer
QNum	Question
10	<p>Consider the following Python statements. Using the “all_users” variable ONLY, print out the surname of the 2nd user.</p> <pre>user_info = [{"name": "macy", "surname": "gray"}, {"name": "john", "surname": "smith"}] all_users = {"users": user_info, "posts": []}</pre>
	Answer

3. The Reddit API

Setting up a developer account

Reddit is a social news and online discussion platform. Reddit consists of thousands of different subreddits. Each subreddit corresponds to a community based around a particular topic or theme. Social media data posted to Reddit can be accessed through the Reddit API.

Access to the Reddit API is achieved in **THREE** steps.

1. Sign up for a normal Reddit account and login.
2. Read the API terms and register your account (<https://www.reddit.com/wiki/api/>)
3. Create an application by completing the form (<https://www.reddit.com/prefs/apps>) NB. this will generate your secret key and APP_ID. Keep these for later.

create application

Please read the API usage guidelines before creating your application. After creating, you will be required to register for production API use.

name

web app A web based application
 installable app An app intended for installation, such as on a mobile phone
 script Script for personal use. Will only have access to the developers accounts

description

about url

redirect uri

create app

Figure 1 The create application form

developed applications



Figure 2 Applications associated with your account

The Python Reddit API Wrapper (PRAW)

[PRAW](#) is an external Python package that simplifies working with the Reddit API in Python applications.

To use external packages in Python they should first be installed. In google collab they can be added to your environment by using the pip command.

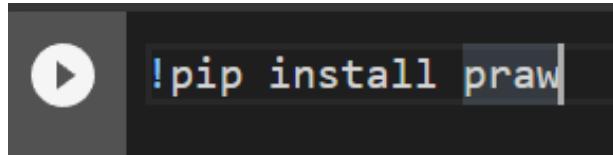
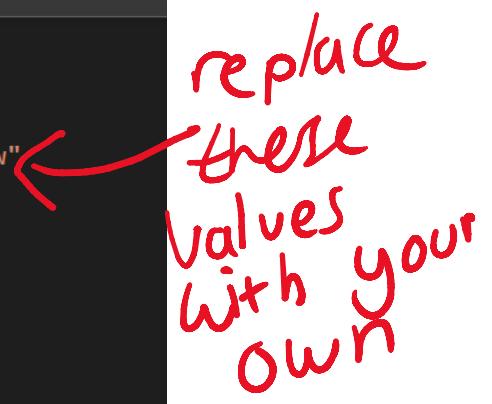


Figure 3 Installing the PRAW package in Google Collab

To use PRAW we first initialise a new praw.Reddit object. This object needs to be initialised with your own APP_ID and secret phrase.



```
import praw
import pprint

SECRET="X0rIaQ7gOera_Vy3RzP796hf2AMhBw"
APP_ID="1Vns2R2WSLPGqXnIg5Dltw"

reddit = praw.Reddit(
    client_id=APP_ID,
    client_secret=SECRET,
    user_agent="Comment Extraction"
)
```

Figure 4 Example setting up the praw.Reddit object

Next, we create a new submission object to represent the post we would like to extract comments from. The submission object requires the full URI of a post on Reddit.

```
url = "https://www.reddit.com/r/britishproblems/comments/11mn8wa/escooters_and_ebikes_could_be_an_amazing_solution/"
submission = reddit.submission(url=url)
```

Figure 5 Setting up the submission object.

To access the top-level comments associated with a submission, we access its comments property.



```
for top_level_comment in submission.comments:
    print(type(top_level_comment))
```

Figure 6 Using a loop to access top-level comments.

here we are
only printing the objects type

When running the above program you should observe that the object type of each comment is itself of type `<class 'praw.models.reddit.comment.Comment'>`.

By reviewing the documentation for the PRAW package we can identify all the properties associated with an individual Comment object. For example, `comment.body` would give us the text associated with a comment (in [Markdown](#) format)

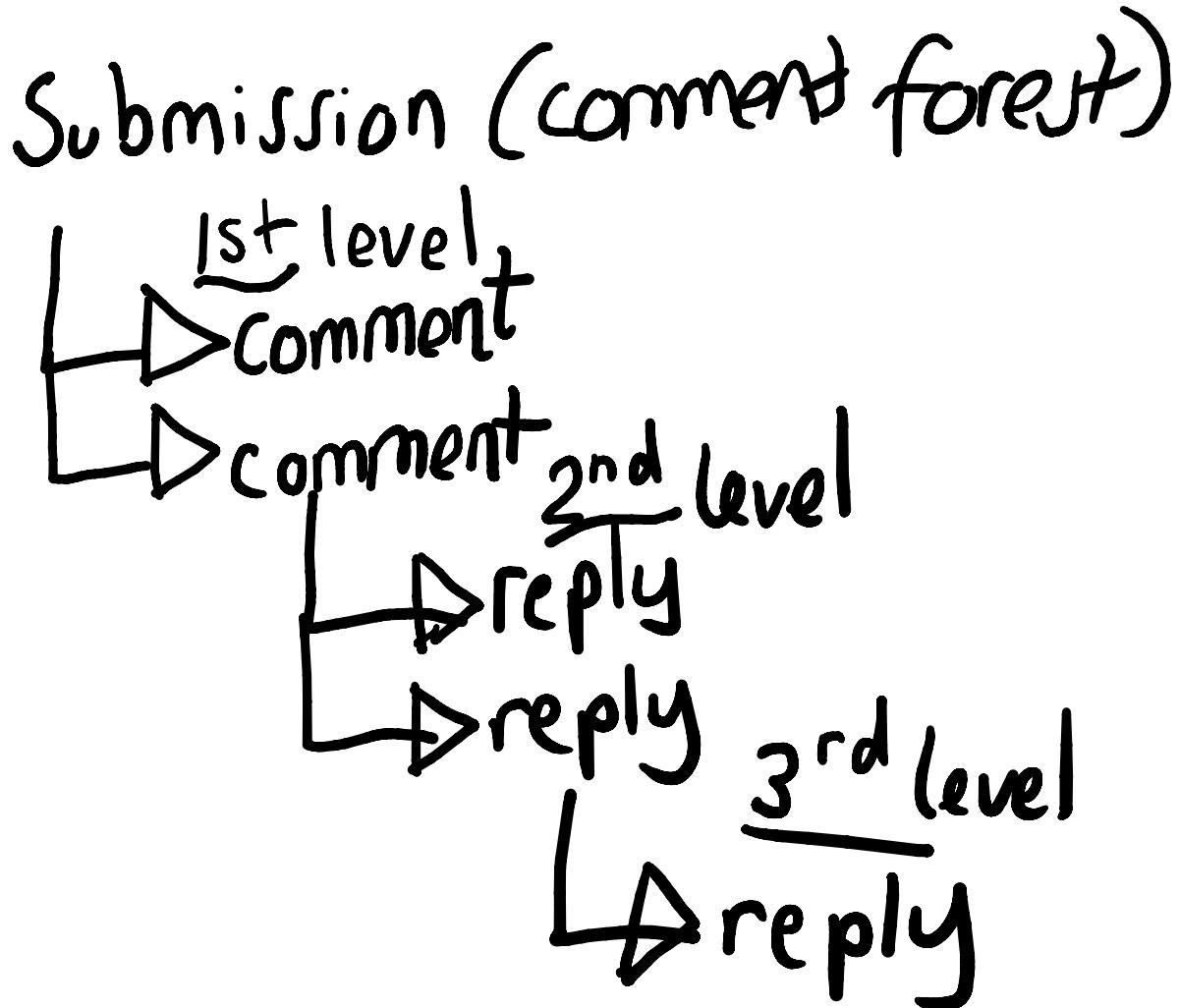
Attribute	Description
<code>author</code>	Provides an instance of <code>Redditor</code> .
<code>body</code>	The body of the comment, as Markdown.
<code>body_html</code>	The body of the comment, as HTML.
<code>created_utc</code>	Time the comment was created, represented in Unix Time .
<code>distinguished</code>	Whether or not the comment is distinguished.
<code>edited</code>	Whether or not the comment has been edited.
<code>id</code>	The ID of the comment.
<code>is_submitter</code>	Whether or not the comment author is also the author of the submission.
<code>link_id</code>	The submission ID that the comment belongs to.
<code>parent_id</code>	The ID of the parent comment (prefixed with <code>t1_</code>). If it is a top-level comment, this returns the submission ID instead (prefixed with <code>t3_</code>).
<code>permalink</code>	A permalink for the comment. <code>comment</code> objects from the inbox have a <code>context</code> attribute instead.
<code>replies</code>	Provides an instance of <code>commentForest</code> .
<code>saved</code>	Whether or not the comment is saved.
<code>score</code>	The number of upvotes for the comment.
<code>stickied</code>	Whether or not the comment is stickied.
<code>submission</code>	Provides an instance of <code>Submission</code> . The submission that the comment belongs to.
<code> subreddit</code>	Provides an instance of <code>Subreddit</code> . The subreddit that the comment belongs to.
<code> subreddit_id</code>	The subreddit ID that the comment belongs to.

Figure 7 Attributes associated with each Comment instance.

By default, accessing the comments property will provide us with access to all the 1st level comments associated with a submission.

In practice there may be more than one level of comments, given the ability of users to respond to each other. 2nd level comments are those that are posted in response to 1st level comments. 3rd level comments are those that were posted in response to 2nd level comments and so forth.

Together these comments resemble a tree-like structure (known as a CommentForest in PRAW)



To access comments further down the tree, we need to loop over the replies of comments in proceeding levels. This is illustrated in the following code, which walks down the tree to the 3rd level.

```
for top_level_comment in submission.comments:  
    for second_level_comment in top_level_comment.replies:  
        for third_level_comment in second_level_comment.replies:  
            print(type(third_level_comment))
```

Figure 8 Accessing third level comments in the CommentForest

*this is an
attribute of
a comment*

Alternatively, we can access the list function of the comments attribute associated with each submission to loop over all comments attached to a submission (regardless of their level in the comment forest).

Keep in mind that since this approach makes many calls to the Reddit API you should be careful not abuse this feature or risk your account being blocked.

```
amount = 10
for comment in submission.comments.list():
    print(type(comment))
    amount -= 1
    if amount == 0:
        break
```

Figure 9 Collapsing the Comment Forest

Attributes of a comment may themselves be other objects. For example, the author associated with each comment is returned as a `<class 'praw.models.reddit.Redditor'>`. This can be verified using the code below.

```
for comment in submission.comments.list():
    print(type(comment.author))
    break
```

Figure 10 Accessing the author property.

The author's name is a property of the Redditor object. The Redditor object is itself attached as a property to each comment.

The following code demonstrates how to print out the author names of 10 comments attached to a submission. The [enumerate](#) function is used to support the loop by returning a counter and the next item from the list each time the loop is run.

```
for n, comment in enumerate(submission.comments.list()):
    print(n, comment.author.name)
    if n == 10:
        break
```

Figure 11 Displaying author names for the first 10 comments



Attribute	Description
<code>comment_karma</code>	The comment karma for the <code>Redditor</code> .
<code>comments</code>	Provide an instance of <code>SubListing</code> for comment access.
<code>submissions</code>	Provide an instance of <code>SubListing</code> for submission access.
<code>created_utc</code>	Time the account was created, represented in Unix Time.
<code>has_verified_email</code>	Whether or not the <code>Redditor</code> has verified their email.
<code>icon_img</code>	The url of the Redditors' avatar.
<code>id</code>	The ID of the <code>Redditor</code> .
<code>is_employee</code>	Whether or not the <code>Redditor</code> is a Reddit employee.
<code>is_friend</code>	Whether or not the <code>Redditor</code> is friends with the authenticated user.
<code>is_mod</code>	Whether or not the <code>Redditor</code> mods any subreddits.
<code>is_gold</code>	Whether or not the <code>Redditor</code> has active Reddit Premium status.
<code>is_suspended</code>	Whether or not the <code>Redditor</code> is currently suspended.
<code>link_karma</code>	The link karma for the <code>Redditor</code> .
<code>name</code>	The Redditor's username.
<code> subreddit</code>	If the <code>Redditor</code> has created a user-subreddit, provides a dictionary of additional attributes. See below.
<code> subreddit["banner_img"]</code>	The URL of the user-subreddit banner.
<code> subreddit["name"]</code>	The fullname of the user-subreddit.
<code> subreddit["over_18"]</code>	Whether or not the user-subreddit is NSFW.
<code> subreddit["public_description"]</code>	The public description of the user-subreddit.
<code> subreddit["subscribers"]</code>	The number of users subscribed to the user-subreddit.
<code> subreddit["title"]</code>	The title of the user-subreddit.

Figure 12 Properties of the `Redditor` object type

4. Pandas

[Pandas](#) is an extremely powerful data manipulation and analysis framework for Python. A core object type in Pandas is the pandas.DataFrame.

A **data frame** is roughly equivalent to a Microsoft Excel worksheet, in the sense that it stores a series of data instances (rows) across a set of fields (columns).

Each **row** in a data frame is associated with a unique identifier, this could be the row number or a unique field present in the dataset. In pandas this unique identifier is known as the **index**.

The screenshot shows a Jupyter Notebook cell with the following code:

```
import pandas as pd  
df = pd.DataFrame()  
print(df)
```

The output below the cell shows:

```
Empty DataFrame  
Columns: []  
Index: []
```

Red handwritten notes on the right side of the screenshot explain the code:

- "this is a shorthand" with an arrow pointing to the import statement.
- "So we don't need to write pandas each time" with an arrow pointing to the pd prefix in the DataFrame constructor.

Figure 13 Creating an empty data frame in pandas

Data frames can be created in various ways. One way is to populate a data frame from a list of dictionaries.

The screenshot shows a Jupyter Notebook cell with the following code:

```
import pandas as pd  
  
row1 = {"name": "Mike", "score": 76}  
row2 = {"name": "John", "score": 52}  
row3 = {"name": "Sara", "score": 91}  
  
rows = [row1, row2, row3]  
  
df = pd.DataFrame(rows)  
print(df)
```

The output below the cell shows:

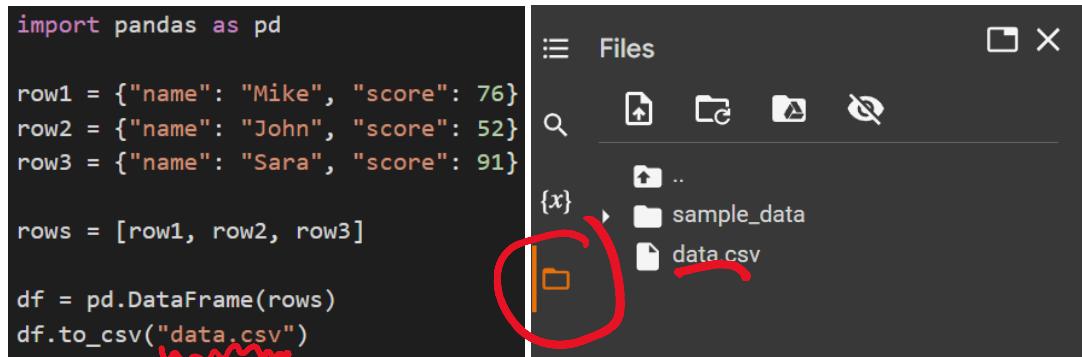
	name	score
0	Mike	76
1	John	52
2	Sara	91

Red handwritten annotations with arrows point to specific parts of the code and output:

- "dictionaries" points to the row definitions (row1, row2, row3).
- "list" points to the rows list.
- "the index" points to the index numbers (0, 1, 2) in the output table.
- "field/column" points to the column names (name, score) in the output table.
- "column value" points to the numerical values (76, 52, 91) in the output table.

Figure 14 Creating a data frame from a list of dictionaries.

Data frames can be saved to CSV (comma separated values) files by calling the `to_csv` method. Similarly, they can also be created by using the `pandas.read_csv(filename)` function.



The screenshot shows a Jupyter Notebook cell containing Python code to create a DataFrame and save it to a CSV file. A red circle highlights the file browser icon in the toolbar, which is used to open the saved file. The file browser shows a folder named 'sample_data' containing a file named 'data.csv'.

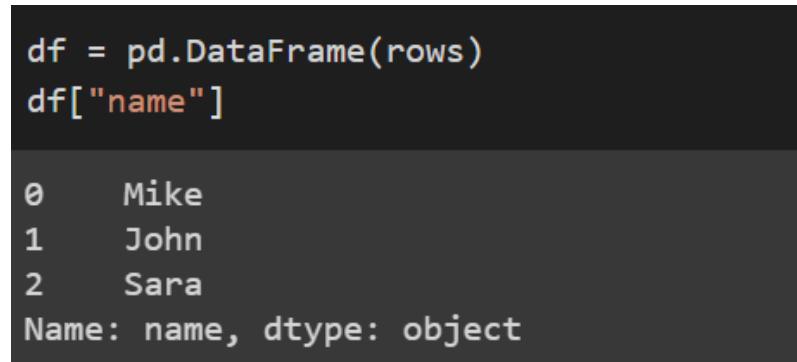
```
import pandas as pd

row1 = {"name": "Mike", "score": 76}
row2 = {"name": "John", "score": 52}
row3 = {"name": "Sara", "score": 91}

rows = [row1, row2, row3]

df = pd.DataFrame(rows)
df.to_csv("data.csv")
```

Individual columns or sets of columns can be accessed using array notation and specifying the column name(s).



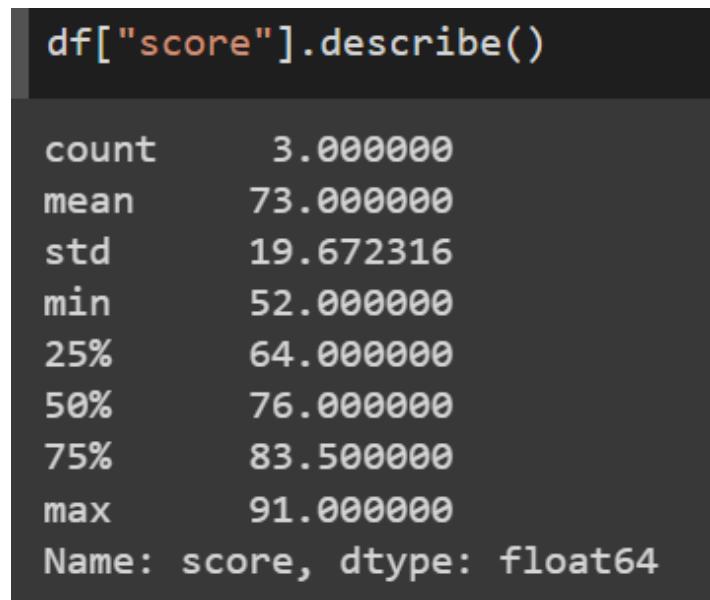
The screenshot shows a Jupyter Notebook cell displaying the selection of a single column from a DataFrame. The code `df["name"]` is run, and the output shows the column data with the header "Name: name, dtype: object".

```
df = pd.DataFrame(rows)
df["name"]

0    Mike
1    John
2    Sara
Name: name, dtype: object
```

Figure 15 Selecting a single column from a data frame.

Various statistics can be computed on a column containing numeric data, e.g. `min()`, `max()`, `mean()`, `std()`, `var()`. The `describe` method returns a set of summary statistics.



The screenshot shows a Jupyter Notebook cell displaying the output of the `describe` method for the 'score' column of a DataFrame. The output provides statistical summary information for the column.

```
df["score"].describe()

count      3.000000
mean      73.000000
std       19.672316
min      52.000000
25%      64.000000
50%      76.000000
75%      83.500000
max      91.000000
Name: score, dtype: float64
```

Adding a new column to a data frame can be done by using array notation and specifying both the name of the new column and its associated value.

```
df = pd.DataFrame(rows)
df["year_of_study"] = 2013
df
```

	name	score	year_of_study
0	Mike	76	2013
1	John	52	2013
2	Sara	91	2013

Figure 16 Adding a new column.

If we want to set the value of a new column according to the values of other fields in the row we can use the apply method. The apply method requires you to specify a function that returns the value of a column given a specific row.

The new field name →

```
df = pd.DataFrame(rows)
df["year_of_study"] = 2023

def score_to_grade(row):
    if row.score >= 70:
        return "distinction"
    if row.score >= 60:
        return "merit"
    if row.score >= 50:
        return "pass"
    return "fail"

df["grade"] = df.apply(score_to_grade, axis=1)
df
```

the function that will determine the value of a field for each row

	name	score	year_of_study	grade
0	Mike	76	2023	distinction
1	John	52	2023	pass
2	Sara	91	2023	distinction

Figure 17 Calculating a new column value based on the value of other columns in the row

Exercises

Using your understanding of the Reddit API and Pandas data frames, attempt the following exercises. It is important that you add comments to your code as you go along.

Q1. Create a Reddit developer account if you haven't already done so.

Select a reddit submission in the reddit.com/r/news subreddit with at least 50 comments.

Note the full URI to this submission. Develop a Python program that prints out the top-level comments for this submission to the console.

Q2. Modify your program so that your program collects 50 comments across *any* level of the Comment Forest. For each comment you should print the body of the comment, the author's name, and the length of the comment in words.

Q3. Modify the program you developed in Q2 so that the comments collected are stored in a data frame. Save this data frame to a CSV file called "comments-dataset.csv". Add a new column to your data frame that is based on the length (in characters) of each comment. Use this new column to determine the min, max and average length of a comment posted.