

CSSE 386 – Data Mining with Programming  
**Rose-Hulman** Institute of Technology

Worksheet 02

Name (Print): \_\_\_\_\_ Section: \_\_\_\_\_

1. For each definition provide the term. Also indicate whether it is a column, row, or table.

Term	Definition
	a characteristic that can be measured or observed
	a collection of information
	data points or observations

2. Define briefly 3 Vs of Big Data

- Volume:
- Variety:
- Velocity:

3. Place lifecycle steps in order: explore, interpret, model, clean, gather

1.                      → 2.                      → 3.                      → 4.                      → 5.

4. Complete the following

```

1  # API from kaggle.....
2  import _____
3  # List the files in the downloaded directory
4  files = os.listdir(path)
5  print("Files in dataset:", _____)
6  csv_file = [file for file in files if file.endswith("_____")]
7  file_path = os.path.join(path, csv_file)
8  df = pd._____ (file_path)
9  print(df._____ ())
10
11 # Load a local file
12 from _____ import files
13 uploaded = files._____ ()
14 df2 = pd._____
15 print(_____.head())

```