

DataEng: Data Validation Activity

Make a copy of this document and use it to record your results. Store a PDF copy of the document in your git repository along with any needed code before submitting for this week.

High quality data is crucial for any data project. This week you'll gain some experience and knowledge of analyzing data sets for quality.

The data set for this week is [a listing of all Oregon automobile crashes on the Mt. Hood Hwy \(Highway 26\) during 2019](#). This data is provided by the [Oregon Department of Transportation](#) and is part of a [larger data set](#) that is often utilized for studies of roads, traffic and safety.

Here is the available documentation for this data: [description of columns](#), [Oregon Crash Data Coding Manual](#)

Data validation is usually an iterative three-step process. First (part A) you develop assertions about your data as a way to make your assumptions explicit. Second (part B) you write code to evaluate the assertions and test the assumptions. This helps you to refine your existing assertions (part C) before starting the whole process over again by creating new assertions (part A again).

Submit: [In-class Activity Submission Form](#)

A. Create Assertions

Access the crash data, review the associated documentation of the data (ignore the data itself for now). Based on the documentation, create English language assertions for various properties of the data. No need to be exhaustive for this assignment, two or more assertions in each category are enough.

1. Create 2+ *existence* assertions. Example, "Every record has a date field".
 - a. Every crash has a date associated with it
 - b. Every crash has at least 1 vehicle associated with it
2. Create 2+ *limit* assertions. The values of most numeric fields should fall within a valid range. Example: "the date field should be between 1/1/2019 and 12/31/2019 inclusive"
 - a. Crash date is in the past, not future.
 - b. Crash participants age is not negative
3. Create 2+ *intra-record check* assertions.
 - a. *There is at least 1 fatality or injury for a given crash*
 - b. Crash month, day and year for each crash level record is a valid date

4. Create 2+ *inter-record check* assertions.
5. Create 2+ *summary* assertions. Example: “every crash has a unique ID”
 - a. Every crash ID is unique
 - b. Most crashes happen outside of school zones
6. Create 2+ *referential integrity* assertions. Example “every crash participant has a Crash ID of a known crash”
 - a. every crash participant has a Crash ID of a known crash
 - b. Every vehicle has a crash id of known crash
7. Create 2+ *statistical distribution assertions*. Example: “crashes are evenly/uniformly distributed throughout the year.”
 - a. Most crashes involve at-least two vehicles
 - b. Most collisions happen at an angle

B. Validate the Assertions

1. Now study the data in an editor or browser. If you are anything like me you will be surprised with what you find. The Oregon DOT made a mess with their data!
2. Write python code to read in the test data and parse it into python data structures. You can write your code any way you like, but we suggest that you use pandas’ methods for reading csv files into a pandas Dataframe
3. Write python code to validate each of the assertions that you created in part A. Again, pandas makes it easy to create and execute assertion validation code.
4. If you are like me you’ll find that some of your assertions don’t make sense once you actually understand the structure of the data. So go back and change your assertions if needed to make them sensible.
5. Run your code and note any assertion violations. List the violations here.

→ None of the assertions that I wrote were violated/failed. All of my assertions related code can be found here:

<https://github.com/pkaran57/data-engineering/blob/main/in-class-assignments/assignment-3/src/crash/CrashDataSet.py#L55>

C. Evaluate the Violations

For any assertion violations found in part B, describe how you might resolve the violation. Options might include “revise assumptions/assertions”, “discard the violating row(s)”, “ignore”, “add missing values”, “interpolate”, “use defaults”, etc.

No need to write code to resolve the violations at this point, you will do that in step E.

If you chose to “revise assumptions/assertions” for any of the violations, then briefly explain how you would revise your assertions based on what you learned.

→ None of the assertions that I thought about in Part A and coded as a part of Part B failed.

D. Learn and Iterate

The process of validating data usually gives us a better understanding of any data set. What have you learned about the data set that you did not know at the beginning of the current ABCD iteration?

→ I did not realize how the data in the excel sheet was structured. Specifically, I did not realize that the excel sheet had multilevel data. The first level being crash, the 2nd level being vehicle and the 3rd level being Participant.

Next, iterate through the process again by going back to Step A. Add more assertions in each of the categories before moving to steps B and C again. Go through the full loop twice before moving to step E.

E. Resolve the Violations

For each assertion violation found during the two loops of the process, write python code to resolve the assertions. This might include dropping rows, dropping columns, adding default values, modifying values or other operations depending on the nature of the violation.

Note that I realize that this data set is somewhat awkward and that it might be best to “resolve the violations” by restructuring the data into proper tables. However, for this week, I ask that you keep the data in its current overall structure. Later (next week) we will have a chance to separate vehicle data and participant data properly.

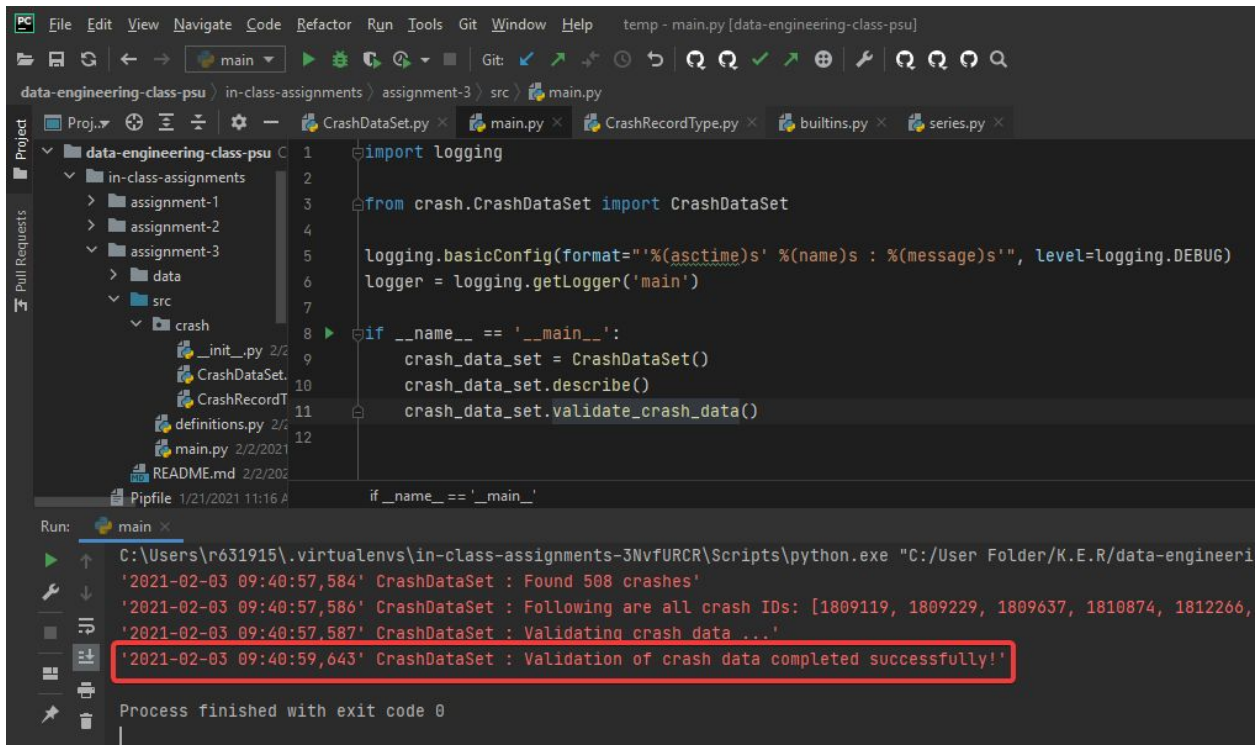
→ I “resolved” some of the data violations by dropping blank age values for instance from the data. Example in code:

<https://github.com/pkaran57/data-engineering/blob/main/in-class-assignments/assignment-3/src/crash/CrashDataSet.py#L73>

E. Retest

After modifying the dataset/stream to resolve the assertion violations you should have produced a new set of data. Run this data through your validation code (Step B) to make sure that it validates cleanly.

→ My validation code validates cleanly:



The screenshot shows a Python IDE with a project named 'data-engineering-class-psu'. The file explorer on the left shows the project structure, including 'in-class-assignments', 'assignment-1', 'assignment-2', 'assignment-3', 'data', 'src', and 'crash'. The main editor displays the code in 'main.py':

```
1 import logging
2
3 from crash.CrashDataSet import CrashDataSet
4
5 logging.basicConfig(format='%(asctime)s' %(name)s : %(message)s', level=logging.DEBUG)
6 logger = logging.getLogger('main')
7
8 if __name__ == '__main__':
9     crash_data_set = CrashDataSet()
10    crash_data_set.describe()
11    crash_data_set.validate_crash_data()
12
```

The terminal output at the bottom shows the execution of the code:

```
Run: main
C:\Users\r631915\.virtualenvs\in-class-assignments-3NvfURCR\Scripts\python.exe "C:/User Folder/K.E.R/data-engineeri
'2021-02-03 09:40:57,584' CrashDataSet : Found 508 crashes'
'2021-02-03 09:40:57,586' CrashDataSet : Following are all crash IDs: [1809119, 1809229, 1809637, 1810874, 1812266,
'2021-02-03 09:40:57,587' CrashDataSet : Validating crash data ...'
'2021-02-03 09:40:59,643' CrashDataSet : Validation of crash data completed successfully!'
Process finished with exit code 0
```

Submit: [In-class Activity Submission Form](#)