

# DataEng S23: Data Transformation

## In-Class Assignment

**Submit:** 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. Submit the in-class activity submission form by Friday at 10:00 pm.

### Initial Discussion Questions

Discuss the following questions among your working group members at the beginning of the week and place your own response into this space. If desired, also include responses from your group members.

1. In the lecture we mentioned the benefits of Data Transformation, but can you think of any problems that might arise with Data Transformation?

Firstly, it may introduce errors if not done carefully. For instance, during data cleaning or normalization, important information could get lost, affecting the accuracy of the results. Secondly, data transformation often requires significant computational resources which might not always be available, especially when dealing with large datasets. Thirdly, the process of transforming data often requires a deep understanding of the data and the problem domain, which might be a challenge if such expertise is not readily available.

2. Should data transformation occur before data validation in your data pipeline or after?

Typically, data validation should occur before data transformation in a data pipeline. The reason is that validation checks for errors, inconsistencies, and anomalies in the raw data. If these are present, they could impact the transformation process, potentially leading to incorrect outputs. Therefore, it's generally advisable to validate your data first, correct any issues identified, and then proceed with data transformation.

### A. Small Sample of TriMet data

Here is sample data for one trip of one TriMet bus on one day (February 15, 2023):

[bc\\_trip259172515\\_230215.csv](#) It's in .csv format not json format, but otherwise, the data is a typical subset of the data that you are using for your class project.

We recommend that you use google Colab or a Jupyter notebook for this assignment, though any python environment should suffice.

Use the [pandas.read\\_csv\(\)](#) method to read the data into a DataFrame.

## B. Filtering

Some of the columns in our TriMet data are not generally useful for our class project. For example, our contact at TriMet told us that the EVENT\_NO\_STOP column is not used and can be safely eliminated for any type of analysis of the data.

Why might we want to filter columns this way instead of using drop()? It's better for efficiency and cleaner code

## C. Decoding

Notice that the timestamp for each breadcrumb record is encoded in an odd way that might make analysis difficult. The breadcrumb timestamps are represented by two columns, OPD\_DATE and ACT\_TIME. OPD\_DATE merely represents the date on which the bus ran, and it should be constant, unchanging for all breadcrumb records for a single day. The ACT\_TIME field indicates an offset, specifically the number of seconds elapsed since midnight on that day.

We're not sure why TriMet represents the breadcrumb timestamps this way. We do know that this encoding of the timestamps makes automated analysis difficult. So your job is to decode TriMet's representation and create a new "TIMESTAMP" column containing a [pandas.Timestamp](#) value for each breadcrumb.

Suggestions:

- Use DataFrame.apply() to apply a function to all rows of your DataFrame
- The applied function should input the two to-be-decoded columns, then it should:
  - create a datetime value from the OPD\_DATE input using datetime.strptime()
  - create a timedelta value from the ACT\_TIME
  - add the timedelta value to the datetime value to produce the resulting timestamp

## D. More Filtering

Now that you have decoded the timestamp you no longer need the OPD\_DATE and ACT\_TIME columns. Delete them from the DataFrame.

## E. Enhance

Create a new column, called SPEED, that is a calculation of meters traveled per second. Calculate SPEED for each breadcrumb using the breadcrumb's METERS and TIMESTAMP values along with the METERS and TIMESTAMP values for the immediately preceding breadcrumb record.

Utilize the [pandas.DataFrame.diff\(\)](#) method for this calculation. diff() allows you to calculate the difference between a cell value and the preceding row's value for that same column. Use diff() to create a new dMETERS column and then again to create a new dTIMESTAMP column. Then use apply() (with a lambda function) to calculate  $SPEED = dMETERS / dTIMESTAMP$ . Finally, drop the unneeded dMETERS And dTIMESTAMP columns.

**Question:** What is the minimum, maximum and average speed for this bus on this trip? (Suggestion: use the Dataframe.describe() method to find these statistics)

## F. Larger Data Set

Here is breadcrumb data for the same bus TriMet for the entire day (February 15, 2023):  
[bc\\_veh4223\\_230215.csv](#)

Do the same transformations (parts B through E) for this larger data set. Be careful, you might need to treat each trip separately. For example, you might need to find all of the unique values for the EVENT\_NO\_TRIP column and then do the transformations separately on each trip.

### Questions:

What was the maximum speed for vehicle #4223 on February 15, 2023?

Where and when did this maximum speed occur?

What was the median speed for this vehicle on this day?

## G. Full Data Set

Here is breadcrumb data for all TriMet vehicles for the entire day (February 15, 2023):  
[bc\\_230215.csv](#)

Do the same transformations (parts B through E) for the entire data set. Again, beware that simple transformations developed in parts B through E probably will need to be modified for the full data set which contains interleaved breadcrumbs from many vehicles.

**Questions:**

What was the maximum speed for any vehicle on February 15, 2023?

Where and when did this maximum speed occur?

Which vehicle had the fastest mean speed for any single trip on this day? Which vehicle and which trip achieved this fastest average speed?