

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?

Information loss: Some changes may result in information loss, particularly when data is reduced in dimensionality. This could have an impact on the analysis of outcomes as well as the overall quality of the decision-making process.

Problems with quality: If the raw data has flaws or inconsistencies, performing transformations may accidentally increase or spread those issues throughout the dataset, potentially resulting in poor-quality outputs.

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

Data transformation and data validation are both crucial steps in a data pipeline, and their sequence is determined by the unique use case and project needs. In general, I feel data validation should come before data transformation.

Before performing transformations, perform data validation to find and repair any flaws or inconsistencies in the raw data. This ensures that the transformed data is based on accurate and reliable input, reducing the risk of propagating errors through the pipeline.

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.

```
import pandas as pd
df = pd.read_csv("bc_trip259172515_230215.csv")
df.head()
```

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.

Use [pandas.DataFrame.drop\(\)](#) to filter the EVENT_NO_STOP column.

For this in-class assignment we won't need the GPS_SATELLITES or GPS_HDOP columns, so drop them as well.

Next, start over and this time try filtering these same columns using the usecols parameter of the read_csv() method.

```
df = pd.read_csv("bc_trip259172515_230215.csv",
usecols=["EVENT_NO_TRIP", "OPD_DATE", "VEHICLE_ID", "METERS", "ACT_TIME",
"GPS_LONGITUDE", "GPS_LATITUDE"])
df.head()
```

Why might we want to filter columns this way instead of using drop()?

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.

CODE in GIT.

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?

The maximum speed of 17.4 meters/second

Where and when did this maximum speed occur?

occurred at latitude 45.505452 and longitude -122.660822.

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

The median speed is 7.2 meters/second.

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?

For vehicle id = 4036

The maximum speed of 32.5 meters/second

Where and when did this maximum speed occur?

occurred at latitude 45.539443 and longitude -122.448413.

The median speed is 8.2 meters/second.

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