

PHW251 Problem Set 7

your name here

today

Part 1

For part 1 of this problem set we will work with motor vehicle crash data from New York City. You can read more about this publicly available data set on their website.

The data file is called “Motor_Vehicle_Collisions_Crashes.csv”. We want you to perform the following:

1. Rename the column names to lower-case and replace spaces with an underscore.
2. Select only:
 - crash_date
 - number_of_persons_injured
 - contributing_factor_vehicle_1
 - vehicle_type_code_1
3. Drop all rows that contain an NA value.
4. Make the values in the vehicle_type_code_1 variable all lowercase and replace the spaces with a dash.
5. Filter the data for vehicles that have a count of at least 500 (appear in the data set 500 times or more)
 - Hints: group_by(), mutate(), n(), filter()
6. Calculate the percentage of accidents by vehicle type
7. Which vehicle group accounted for 1.55% (0.0155) of the accidents?

We have grouped the questions below to push you to perform commands with less code. As you’re building your code we recommend going line by line to test, then combining to perform multiple steps in one command.

Questions 1-3

```
# YOUR CODE HERE
```

Questions 4-5

```
# YOUR CODE HERE
```

Question 6

```
# YOUR CODE HERE
```

Question 7

WRITE YOUR ANSWER HERE

Part 2

For this part we will work with four tables that are relational to each other. The following keys link the tables together:

- patient_id: patients, schedule
- visit_id: schedule, visits
- doctor_id: visits, doctors

Question 8

You've been asked to collect information on patients who are actually on the schedule. To start this task, you need to join the patient data to the schedule data, since we only want to keep the observations that are present in both the patient data AND the schedule data.

Which kind of join do you use?

WRITE YOUR ANSWER HERE

How many observations do you see in your joined data set? Notice that some patients have multiple visits.

YOUR CODE HERE

WRITE YOUR ANSWER HERE

Question 9

In the visits data, we have a variable called “follow_up” where Y means a follow-up is needed and N means a follow-up is not needed. How many patients require a follow-up? You will want to first make a join and then subset. Start with the data frame created in the previous question.

```
# YOUR CODE HERE
```

Which join did you use?

WRITE YOUR ANSWER HERE

How many patients need a follow-up?

WRITE YOUR ANSWER HERE

Question 10

Which doctors do these patients need follow-up with? Print out each doctor's name.

YOUR CODE HERE

Which join did you use?

WRITE YOUR ANSWER HERE