

PHW251 Take-home Midterm Exam

STUDENT COPY

October 2023

Instructions

Where to put your code and answers

This R Markdown document gives the instructions for the take-home portion of your midterm exam. Submit your completed exam by filling in the “responses.rmd” file. When completing this Rmd, please include your answer for each question where it says **ANSWER HERE:**. Please also provide your code in the indicated code chunks. The code portion will be used to assign partial credit.

How to answer questions

Please carefully read the **Necessary Coding Steps** instructions for each question.

All questions below should be answered using R. Unless otherwise specified, you may use any method (base R, tidyverse, or other) to answer these questions. Please type out your answers in the specified area **ANSWER HERE:** (even if the answer is also available in your code chunk). Code will be used to give partial credit for incorrect answers.

Submitting

When you are done, knit the “responses.rmd” file to PDF and load it into Gradescope before the due date. Additionally, keep in mind best practices including ensuring code does not run off the page and not printing entire dataframes within your final PDF. **When knitted, your PDF should be approximately one page per question** (points will be deducted if too long with unnecessary output).

Scoring

Your exam will be worth **15 points total**. Each question indicates the number of points for that question.

Data

For all questions below, use the exam_data_2023_version1.csv file that is saved on DataHub/GitHub at PHW251_Fall2023/midterm/data/exam_data_2023_version1.csv. This is a real dataset from the California Health and Human Services Open Data Portal, but has been altered slightly for the purpose of this exam. The dataset contains counts of emergency department encounters at California medical facilities.

The file includes the following columns:

- Year

- **OSHPD ID**
- **Facility Name**
- **County Name**
- **ER Service Level Desc:** Level of ER service. Options include: BASIC, COMPREHENSIVE, STANDBY, NOT APPLICABLE
- **Type:** Specifies encounter type. Options include: ED_Visit (Encounter in which patient is treated in the Emergency Department and then released), ED_Admit (Encounter in which the patient is initially treated in the Emergency Department and then admitted to the same hospital for continued inpatient care). Categories are mutually exclusive.
- **Count**

Exam Questions

Question 1

Necessary Coding Steps

Import the csv data file. Check what data type each column is.

1A. What are the data types of each column when you read the data into R (numeric, factor, logical, character, etc)? [1 pt]

ANSWER:

- Year: ?
- YOSHPD ID: ?
- Facility Name: ?
- County Name: ?
- ER Service Level Desc: ?
- Type: ?
- Count: ?

Question 2

Necessary Coding Steps

Notice the column names are not reading in in a very user-friendly way. Rename all columns to align with best practices for naming columns (lowercase with underscores in place of spaces).

2A. What are the new column names? Please list the new column names below. [1 pt]

ANSWER: ?

Question 3

Questions 3-7 are designed to build off each other.

Necessary Coding Steps

Using the data frame from question 2, create a new data frame that limits the data frame to only contain rows where the type of service was “basic” and year is between 2014 and 2018 (inclusive of these years).

3A. How many records are in the new subsetting dataset? [1 pt]

ANSWER: ?

Question 4

Necessary Coding Steps

Using the data frame created in question 3, create a new column called `total_encounters` by grouping OSHPD ID and Year and then summing the values in the `count` column to get total encounters. (*Hint: After adding this column your data frame should contain the same number of rows that it had before you added the column.*)

4A. What is the value of `total_encounters` for LOS ALAMITOS MEDICAL CENTER in 2018? [1 pt]

Answer: ?

Necessary Coding Steps

Create another new column called `pct_encounter_type` that calculates the percent of ED encounters that were visits or admits. Display the percentage as multiplied by 100 and rounded to 1 decimal (for example, 35.1% would be displayed as 35.1).

4B. What is the value of `pct_encounter_type` for ED admits at SUTTER DAVIS HOSPITAL in 2016? [1 pt]

Answer: ?

Question 5

Necessary Coding Steps

Using the data frame created in question 4, first create a subset table that only includes rows for ED admits. Then use the arrange function to order the data frame to display rows first by lowest to highest year and then by highest to lowest value of `pct_encounter_type` (at the same time).

5A. Show a single line of code that can be used for this arrange step. [1 pt]

ANSWER: ?

5B. What code would you use to obtain only the facility names for the first 5 rows of the dataset created in question 5A (facilities with highest values in the `pct_encounter_type` column for the first year in the data frame)? [1 pt]

ANSWER: ?

Question 6

Necessary Coding Steps

Using the dataset created in question 5A, find the average (mean) value of percent of encounter types (`pct_encounter_type`) that were admits among all facilities from 2014 to 2018. Use this mean value to create another new column called `above_below_avg` that categorizes facilities with `pct_encounter_type` equal to or above average as “above”; otherwise, categorize as “below”.

6A. What is the average (mean) percentage of encounters that are ED admits for all facilities from 2014 to 2018? [1 pt] ANSWER HERE: ?

6B. What is the value of `above_below_avg` for ER admits at SOUTH COAST GLOBAL MEDICAL CENTER for 2018? ANSWER: ?

Question 7

Necessary Coding Steps

Restrict the dataset created in 6A to only include records for TULARE, CONTRA COSTA, and SANTA CLARA facilities for the year 2016.

7A. How many records remain? [1 pt] ANSWER: ?

7B. Using the data frame created in question 7A, what hospital has the highest percent of encounters that are admits? [1 pt] ANSWER: ?

Question 8

Necessary Coding Steps

For questions 8 & 9, please use the data frame from the end of question 2. To start, create a new subset data frame that only includes records for encounters that were ED visits (not admits) in the year 2018. Additionally, only include the following columns: Facility Name, County, ER Service Level visits, Type, and Count.

Create a new column called `county_visit_total` that contains the total number of ED visits for each county. (*Hint: this data frame should contain 1 row per county.*) Re-order the table to display the county with the highest number of ED visits at the top of the table.

8A What county has the 10th highest total number of ED visits in 2018? [1 pt] ANSWER HERE: ?

Question 9

Necessary Coding Steps

Building on to the data frame created in question 8, create a new column called `visit_category` indicating the categorical level of ER visits utilization (High, Medium, Low, Very Low) in each county. The categories should be defined as:

- “High”: > 178649 ED visits
- “Medium”: > 66521 ED visits
- “Low”: > 22026 ED visits
- “Very Low”: ≤ 22026 ED visits

Create a final table that summarizes the number of counties in each category (*Hint: this table should only have 4 rows*).

9A. How many counties are in the “low” coverage category? [1 pt]

ANSWER HERE: ?

Question 10

Necessary Coding Steps

For question 10, please use the data frame from the end of question 2 – NOT the data frames that you created in questions 8 and 9.

Create a subset dataset with only records for 2013 and basic ER service level. Keep only the following columns: Facility Name, County Name, Type, and Count. Pivot the dataset to create columns for each of the ED encounter types; these columns should each contain the counts of encounters for each type (admit and visit).

10A. How many records are in the dataset after the pivot? [1 pt] ANSWER: ?

10B. Include the line of code used to perform the pivot. Make sure to include the function name as well as the arguments used. [1 pt] ANSWER HERE: ?

EXTRA CREDIT

[2 points total]

Question 11 (extra credit)

Complete questions 8 and 9 using only one dplyr call. In other words, start with the data frame from the end of question 2, perform the necessary subsetting, grouping, and summarizing with the end goal of producing a table that displays the number of counties in each visit count category.

Please include sufficient code for the teaching team to be able to run your code, if needed. This means either including the import statement for the csv before the dplyr call, or including the import statement as part of your dplyr call.

Hint: Including more than one `group_by()` in a single call may also require the use of `ungroup()`.

Paste the single dplyr call below. [1 pt]

Question 12 (extra credit)

Include code that uses the `kable` package to print the final table for question 9 in a print-friendly format (easy to read with meaningful column names and rows in descending order from High to Very Low). [1 pt]