

Homework Assignment 3

STA 141A

Due Tuesday, November 14 by 5:00 pm

Description

Bike sharing programs have become popular in big cities over the last decade. In the most common arrangement, customers can rent a bike from an automated "dock" and return it later to any other dock in the city. The charge the customer is based on how long they used the bike.

The largest bike sharing program in the United States is in New York City, but bike shares have also recently formed in Los Angeles and the San Francisco Bay Area. The bike sharing programs in many cities publish anonymized data about trips and stations.

In this assignment, you will analyze open data provided by the Los Angeles and Bay Area bike shares. The data sets are provided in a ZIP compressed file. The ZIP file contains:

- 1 CSV file with Bay Area bike share stations.
- 1 CSV file with Bay Area bike share trips.
- 1 CSV file with Los Angeles (Metro) bike share stations.
- 5 CSV files with Los Angeles (Metro) bike share trips.

Questions

Use R to find answers to all of the following questions (that is, don't do any by hand or by point-and-click). Save your code in an R script.

1. Write a function that loads the Bay Area bike share trip data from a CSV file, converts the columns to appropriate data types, and then saves the tidied data frame to an RDS file. Your function should have arguments to set the path for the input CSV file and the output RDS file.

Write a second function that does the same thing for the Bay Area bike share station data.

2. Create a map that shows the locations of the Bay Area bike share stations in San Francisco (only). Label each station with its name. Make the size of each point correspond to the number of trips started from that station. Discuss what you can conclude from the map.
3. Write a function that loads the Los Angeles bike share trip data from the 5 provided CSV files, binds them into one data frame, converts the columns to appropriate data types, and saves the tidied data frame to an RDS file. Your function should have arguments to set the path for the input directory and the output RDS file. Keep your function short and simple by using an apply function rather than repeating code.

Write a second function that loads, tidies, and saves the Los Angeles bike share station data.

4. Create a map that shows the locations of the Los Angeles bike share stations near downtown Los Angeles (only). Label each station with its name. Make the size of each point correspond to the number of trips started from that station. Discuss what you can conclude from the map.

5. How do trip frequency, distance, and duration change at different times of day? Investigate for both the Bay Area bike share and the Los Angeles bike share. Compare your findings. The `geosphere::distGeo()`¹ function can compute distances for longitude and latitude coordinates.
6. For Bay Area bike share trips in San Francisco, how does bearing (angle) change at different times of day? What can you conclude about traffic patterns in the city? The `geosphere::bearing()` function can compute bearings for longitude and latitude coordinates.

Assemble your answers into a report. Please do not include any raw R output. Instead, present your results as neatly formatted² tables or graphics, and write something about each one. You must **cite your sources**. Your report should be **no more than 5 pages** including graphics, but excluding code and citations.

What To Submit

Submit a digital copy on Canvas. The digital copy must contain your report (as a PDF) and your code (as one or more R scripts).

Additionally, submit a printed copy to the box in the statistics department office³. The printed copy must contain your report and your code (in an appendix). Please print double-sided to save trees. It is your responsibility to make sure the graphics are legible in the printed copy!

Data Documentation

The features in the Bay Area bike share trips data set are:

<code>trip_id</code>	unique id of trip
<code>duration_sec</code>	duration of trip (in seconds)
<code>start_date</code>	time trip started (in PST)
<code>start_station_name</code>	name of station where trip started
<code>start_station_id</code>	unique id of station where trip started
<code>end_date</code>	time trip ended (in PST)
<code>end_station_name</code>	name of station where trip ended
<code>end_station_id</code>	unique id of station where trip ended
<code>bike_number</code>	unique id of bike used
<code>zip_code</code>	home zip code of rider (potentially unreliable)
<code>subscriber_type</code>	Subscriber = annual or 30-day member; Customer = 24-hour or 3-day member

The features in the Bay Area bike share stations data set are:

<code>station_id</code>	unique id of station
<code>name</code>	name of station
<code>latitude</code>	latitude of station
<code>longitude</code>	longitude of station
<code>dockcount</code>	total number of docks at the station
<code>landmark</code>	city
<code>installation_date</code>	date the station was installed

A station row may be partially duplicated in the data set if the station was moved.

The features in the Metro bike share data sets are documented here: <https://bikeshare.metro.net/about/data/>

¹The `::` is R notation for a function in a specific package. Read it as `PACKAGE::FUNCTION()`.

²See the graphics checklist on Canvas.

³4th floor of Mathematical Sciences Building

Relevant Functions

All of the functions from previous assignments, as well as:

`cut()`, `saveRDS()`, `strptime()`, `as.data.frame()`, `lapply()`

The functions in the packages *ggmap*, *geosphere*, *ggplot2*, *lubridate*, and *readr* may also be useful, although you are not required to use these packages.