

Data Wrangling (Data Preprocessing)

Practical assessment 2

s9001731 Mark randall

2024-05-22

Assignment cover sheet

for use when submitting work for assessment

School

Master of Data Science

Program name

MC267

Course/unit name

Data Wrangling(2410)

Course/unit code

MATH2349

National Unit of Competency (UOC) title (VE only)

National Unit of Competency (UOC) code (VE only)

Name of lecturer/teacher

Dr. Sona Tehari

Name of tutor/marker/assessor

Teaching Team -

Assignment no.

3

Due date (DD/MM/YYYY)

31/05/2024

Class daytime

Mon - Wed

Campus

City-Online

Student/s

Randall

Mark

s9001731

Family name

Given name

Student no.

Family name

Given name

Student no.

Family name

Given name

Student no.

Family name

Given name

Student no.

Family name

Given name

Student no.

Family name

Given name

Student no.

Declaration and statement of authorship

1. I/we have not impersonated, or allowed myself to be impersonated by, any person for the purposes of this assessment.

2. This assessment is my/our original work and no part of it has been copied from any other source except where due acknowledgement is made.

3. No part of this assessment has been written for me/us by any other person except where such collaboration has been authorised by the lecturer/teacher concerned.

4. Where this work is being submitted for individual assessment, I declare that it is my original work and that no part has been contributed by, produced by or in conjunction with another student.

5. I/we give permission for my assessment response to be reproduced, communicated compared and archived for the purposes of detecting plagiarism.

6. I/we give permission for a copy of my assessment to be retained by the university for review and comparison, including review by external examiners.

7. I/we understand that:

Plagiarism is the presentation of the work, idea or creation of another person as though it is your own. It is a form of cheating and is a very serious academic offence that may lead to exclusion from the University. Plagiarised material can be drawn from, and presented in, written, graphic and visual form, including electronic data and oral presentations. Plagiarism occurs when the origin of the material used is not appropriately cited.

Plagiarism includes the act of assisting or allowing another person to plagiarise or to copy my work.

I/we agree and acknowledge that:

I/we have read and understood the Declaration and Statement of Authorship above.

If I/we do not agree to the Declaration and Statement of Authorship in this context and a signature is not included below, the assessment outcome is not valid for assessment purposes and will not be included in my final result for this course.

Visit the student essential [Academic Integrity](#) page for more information on academic integrity at RMIT, including your responsibilities, breaches and penalties for academic misconduct.

Student signature/s required on page 2.

Assignment cover sheet

for use when submitting work for assessment

Student signature/s

I/we declare that I/we have read and understood the declaration and statement of authorship.

1

Digitally signed by Mark Randall

Date: 2024.05.16 21:24:45 +1000

2

3

4

5

6

Further information relating to academic integrity breaches, and penalties that range from a notation on your student record to expulsion from the University, is contained in the [Academic Integrity Policy](#), [Student Conduct Policy – Schedule 1](#), and the [Student Conduct Regulations](#).

Integrity | Assignment Cover Sheet for use when submitting work for assessment

0984 1221 | 1 of 2

Integrity | Assignment Cover Sheet for use when submitting work for assessment

0984 1221 | 2 of 2

Student names, numbers and percentage of contributions

Table 1: Group information

Student name	Student number	Percentage of contribution
Mark Randall	s9001731	100

1

## Library Load

Package 1 : OpenStreetMap (Fellows & JMapView library by Jan Peter Stotz 2023)  
Package 2 : tidyterra (Hernangómez 2024)  
Package 3 : maptiles (Giraud 2024)  
Package 4 : sf (Pebesma 2024)  
Package 5 : sn (Azzalini 2023)  
Package 6 : stats4 (R Core Team 2024a)  
Package 7 : moments (Komsta & Novomestky 2022)  
Package 8 : ggnewscale (Campitelli 2024)  
Package 9 : Hmisc (Harrell Jr 2024)  
Package 10 : validate (van der Loo & de Jonge 2024)  
Package 11 : deducorrect (van der Loo, de Jonge & Scholtus 2015)  
Package 12 : editrules (de Jonge & van der Loo 2024)  
Package 13 : igraph (Csárdi et al. 2024)  
Package 14 : deductive (van der Loo & de Jonge 2021)  
Package 15 : tidysselect (Henry & Wickham 2024)  
Package 16 : rvest (Wickham 2024)  
Package 17 : here (Müller 2020)  
Package 18 : glue (Hester & Bryan 2024)  
Package 19 : magrittr (Bache & Wickham 2022)  
Package 20 : lubridate (Spinu, Grolemund & Wickham 2023)  
Package 21 : forcats (Wickham 2023a)  
Package 22 : stringr (Wickham 2023b)  
Package 23 : purrr (Wickham & Henry 2023)  
Package 24 : tibble (Müller & Wickham 2023)  
Package 25 : ggplot2 (Wickham et al. 2024)  
Package 26 : tidyverse (Wickham 2023c)  
Package 27 : kableExtra (Zhu 2024)  
Package 28 : knitr (Xie 2024)  
Package 29 : readxl (Wickham & Bryan 2023)  
Package 30 : readr (Wickham, Hester & Bryan 2024)  
Package 31 : dplyr (Wickham et al. 2023)  
Package 32 : tidyr (Wickham, Vaughan & Girlich 2024)  
Package 33 : openxlsx (Schauburger & Walker 2023)  
Package 34 : stats (R Core Team 2024b)  
Package 35 : graphics (R Core Team 2024c)  
Package 36 : grDevices (R Core Team 2024d)  
Package 37 : utils (R Core Team 2024e)  
Package 38 : datasets (R Core Team 2024f)  
Package 39 : methods (R Core Team 2024g)  
Package 40 : base (R Core Team 2024h)

## Abstract

*“Most vehicular accidents in Victoria involve a male driver between 18 to 30 years of Age and a high powered car.”*

This project will use some empirical data collected by the Victorian State Government(Vic 2024) to examine

this statement.

## Executive Summary

The Victoria Road Crash Data URL(Vic 2024) contains nine comma-separated value (csv) and one geo spatial java script object notation (GeoJSON) file. These were downloaded to a Data folder for examination. The files are:

```
csvFileNames <- list.files("../Data",pattern = "*.csv", full.names = TRUE)
fullFileNames <- list.files("../Data", full.names = TRUE)
for (file in fullFileNames) {
  cat(paste("-\t", file,"\n"))
}
```

- ../Data/ACCIDENT.csv
- ../Data/ACCIDENT\_EVENT.csv
- ../Data/ACCIDENT\_LOCATION.csv
- ../Data/ATMOSPHERIC\_COND.csv
- ../Data/NODE.csv
- ../Data/Order\_9WAPHM.zip
- ../Data/PERSON.csv
- ../Data/ROAD\_SURFACE\_COND.csv
- ../Data/SUB\_DCA.csv
- ../Data/VEHICLE.csv
- ../Data/VICTORIAN\_ROAD\_CRASH\_DATA.geojson

The URL indicates that the metadata was updated 29 April 2024, data observations updated as at 27 November 2024 and that observations temporal start was 1 January 2012.

## Data

Provide explanations here.

```
# Import the data, provide your R codes here.
```

## Understand

```
# This is the R chunk for the Understand Section
```

Provide explanations here.

## Tidy & Manipulate Data I

```
# This is the R chunk for the Tidy & Manipulate Data I
```

Provide explanations here.

## Tidy & Manipulate Data II

```
# This is the R chunk for the Tidy & Manipulate Data II
```

Provide explanations here.

## Scan I

```
# This is the R chunk for the Scan I
```

Provide explanations here.

## Scan II

```
# This is the R chunk for the Scan II
```

Provide explanations here.

## Transform

```
# This is the R chunk for the Transform Section
```

Provide explanations here.

## Bibliography

Azzalini, A 2023, *Sn: The skew-normal and related distributions such as the skew-t and the SUN*,.

Bache, SM & Wickham, H 2022, *Magrittr: A forward-pipe operator for r*,.

Campitelli, E 2024, *Ggnewscale: Multiple fill and colour scales in 'ggplot2'*,.

Csárdi, G, Nepusz, T, Traag, V, Horvát, S, Zanini, F, Noom, D & Müller, K 2024, *Igraph: Network analysis and visualization*,.

de Jonge, E & van der Loo, M 2024, *Editrules: Parsing, applying, and manipulating data cleaning rules*,.

Fellows, I & JMapView library by Jan Peter Stotz, using the 2023, *OpenStreetMap: Access to open street map raster images*,.

Giraud, T 2024, *Maptiles: Download and display map tiles*,.

Harrell Jr, FE 2024, *Hmisc: Harrell miscellaneous*,.

Henry, L & Wickham, H 2024, *Tidysselect: Select from a set of strings*,.

Hernangómez, D 2024, *Tidyterra: 'Tidyverse' methods and 'ggplot2' helpers for 'terra' objects*,.

Hester, J & Bryan, J 2024, *Glue: Interpreted string literals*,.

Komsta, L & Novomestky, F 2022, *Moments: Moments, cumulants, skewness, kurtosis and related tests*,.

Müller, K 2020, *Here: A simpler way to find your files*,.

Müller, K & Wickham, H 2023, *Tibble: Simple data frames*,.

Pebesma, E 2024, *Sf: Simple features for r*,.

R Core Team 2024a, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

R Core Team 2024b, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

R Core Team 2024c, *R: A language and environment for statistical computing*, R Foundation

for Statistical Computing, Vienna, Austria.

R Core Team 2024d, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

R Core Team 2024e, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

R Core Team 2024f, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

R Core Team 2024g, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

R Core Team 2024h, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria.

Schauberger, P & Walker, A 2023, *Openxlsx: Read, write and edit xlsx files*,.

Spinu, V, Grolemund, G & Wickham, H 2023, *Lubridate: Make dealing with dates a little easier*,.

van der Loo, M & de Jonge, E 2021, *Deductive: Data correction and imputation using deductive methods*,.

van der Loo, M & de Jonge, E 2024, *Validate: Data validation infrastructure*,.

van der Loo, M, de Jonge, E & Scholtus, S 2015, *Deducorrect: Deductive correction, deductive imputation, and deterministic correction*,.

Vic, D 2024, 'Victoria road crash data', vol. 2024.

Wickham, H 2023a, *Forcats: Tools for working with categorical variables (factors)*,.

Wickham, H 2023b, *Stringr: Simple, consistent wrappers for common string operations*,.

Wickham, H 2023c, *Tidyverse: Easily install and load the 'tidyverse'*,.

Wickham, H 2024, *Rvest: Easily harvest (scrape) web pages*,.

Wickham, H & Bryan, J 2023, *Readxl: Read excel files*,.

Wickham, H, Chang, W, Henry, L, Pedersen, TL, Takahashi, K, Wilke, C, Woo, K, Yutani, H, Dunnington, D & van den Brand, T 2024, *ggplot2: Create elegant data visualisations using the grammar of graphics*,.

Wickham, H, François, R, Henry, L, Müller, K & Vaughan, D 2023, *Dplyr: A grammar of data manipulation*,.

Wickham, H & Henry, L 2023, *Purrr: Functional programming tools*,.

Wickham, H, Hester, J & Bryan, J 2024, *Readr: Read rectangular text data*,.

Wickham, H, Vaughan, D & Girlich, M 2024, *Tidyr: Tidy messy data*,.

Xie, Y 2024, *Knitr: A general-purpose package for dynamic report generation in r*,.

Zhu, H 2024, *kableExtra: Construct complex table with 'kable' and pipe syntax*,.