

# Cleaning Datasets for Analysis

## Contents

<b>A brief explanation of messy data</b>	<b>1</b>
Reproducible data cleaning workflows . . . . .	2
<b>Load data from file</b>	<b>2</b>
Comma separated values ( <b>.csv</b> ) . . . . .	3
Excel spreadsheet ( <b>.xlsx</b> ) . . . . .	3
Single tab . . . . .	3
Multiple tabs . . . . .	3
<b>Common messes and their fixes</b>	<b>3</b>
Missing values . . . . .	3
Interpolate . . . . .	3
Treat as zero . . . . .	3
Ignore . . . . .	3
Inconsistent types . . . . .	3
Detecting type . . . . .	3
Coercing type . . . . .	3
Inconsistent formats . . . . .	3
Detecting type . . . . .	3
Correcting type . . . . .	3
Typos . . . . .	3
Off-by-one errors . . . . .	3
Extra white space . . . . .	3

## A brief explanation of messy data

“Messy data” are data that do not conform to the expected formats required to be read by a machine, or have sensible but inconsistent formats that complicate machine reading. For example, you might recognize “2022-02-03” and “2/3/2022” as two different date formats for February 3, 2022. But a computer program would need explicit instruction on date formats to understand that these two character strings share the same interpretation.

## Reproducible data cleaning workflows

Scripted data cleaning workflows allow you to document how you approach data cleaning in a reproducible and easy-to-share manner. Data cleaning scripts are self-documenting; if you save your R scripts for data cleaning then you have a record of your approach, the decisions you made, and a means of reproducing your cleaning steps. While it may seem easier to clean messy data manually in a spreadsheet application, the reproducibility and transparency benefit of scripted workflows makes the scripted approach the widely preferred method for data analytics.

## Load data from file

Name	Birthday
Harry Potter	31st July, 1980
Ronald Weasley	1st March, 1980
Hermione Granger	19 September, 1979
Neville Longbottom	30th July, 1980
Ginny Weasley	11th August, 1981
Luna Lovegood	13th February, 1981
Fred and George Weasley	1st April, 1978
James Potter	27th March, 1960
Lily Potter	30th January, 1960
Sirius Black	3rd November, 1959
Remus Lupin	10th March, 1960
Minerva McGonagall	4th October
Rubeus Hagrid	6th December
Severus Snape	9th January, 1960
Voldemort	31st December

Comma separated values (.csv)

Excel spreadsheet (.xlsx)

Single tab

Multiple tabs

## Common messes and their fixes

Missing values

Interpolate

Treat as zero

Ignore

Inconsistent types

Detecting type

Coercing type

Inconsistent formats

Detecting type

Correcting type

Typos

Off-by-one errors

Extra white space