

# Data Appendix to “Name of your paper”

Avery Hammond

## Contents

<b>1</b>	<b>Appendix description</b>	<b>1</b>
<b>2</b>	<b>Instructions for Use</b>	<b>1</b>
<b>3</b>	<b>Raw data</b>	<b>2</b>
3.1	Dataset description . . . . .	2
<b>4</b>	<b>Data Processing and Combination</b>	<b>12</b>
<b>5</b>	<b>Analysis Variables</b>	<b>12</b>
<b>6</b>	<b>Discussion of Data</b>	<b>12</b>

```
## Warning in fun(libname, pkgname): couldn't connect to display ":0"
```

## 1 Appendix description

*Your Data Appendix should begin with a brief statement explaining its purpose like the following one.*

This Data Appendix documents the data used in “Papter Title”. It was prepared in a Rmarkdown document that contains both the documentation and the R code used to prepare the data used in the final estimation. It also includes descriptive statistics for both the original data and the final dataset, with a discussion of any issues of note.

The datasets used directly by the final analysis are saved in **processed-data/** at the end of this file.

*Note: this document structure will require you to re-run steps of your analysis multiple times. If your code takes a long time, please come talk with me about strategies to reduce run time or save earlier results.*

## 2 Instructions for Use

This document includes instructions for how to create your Data Appendix. Outside of this section, instruction paragraphs are listed in *italics* (like the first paragraph above). Instructions should be removed before submission.

To start creating your own data appendix, follow these steps:

1. Replace the title and author in the section at the top of the file (called the YAML).

2. Commit your changes with a message like “customizing data appendix”.
3. Delete this instruction section of the document.
4. Remove any other instructions in italics and examples from the completed sections of the document.

Remember that you will submit your assignment by committing and then pushing your versions to your repository. I encourage you to commit your changes often as you work, but there are three specific points at which you need to both submit and push changes, corresponding to course deadlines:

1. You must submit a version with the original data section completed by the Data Appendix 1 deadline. This will include the .Rmd file, the .pdf file, and the html data summary files stored in the output folder.
2. You must submit a version with all parts completed by the Data Appendix 2 deadline.
3. You must submit a final version of this document that is consistent with your final paper by the final project deadline.

While creating your data appendix, refer regularly to the assignment descriptions posted on Moodle.

A few tips:

- When creating a list like this one, be sure to put an empty line above the list. If you don’t do this, your entries won’t be formatted a list.
- Make sure you have empty lines above and below section and subsection headings.
- When creating numbered lists, you can number all items in your list with 1. Rmarkdown will number them sequentially when it creates your final document.

## 3 Raw data

*Each dataset you use will have its own documentation section. The next subsection in this document (Dataset description) is a template. You can copy this section and paste it into your document each time you need to add a section for a new dataset. Note that each line in the Dataset description section **must** end with two spaces. This section documents the datasets used in this analysis.*

### 3.1 Dataset description

**Citation:** Put citation here in APA or other consistent format that you will use throughout the project. Include a hyperlink if applicable.

**DOI:** If the dataset has a documentation identified (DOI) assigned put it here.

**Date Downloaded:** Identify when you downloaded the dataset.

**Filename(s):** raw\_data/filename.csv *If you have a large number of files you can use a patten (see visit data below)* **Unit of observation:** What distinguishes different rows in your dataset?

**Dates covered:** What time frame does the data cover?

#### 3.1.1 To obtain a copy

Describe in a step-by-step fashion how an interested user could obtain the data.

#### 3.1.2 Importable version (if necessary)

**Filename(s):** importable-data/filename-importable.csv

In some cases the raw data is not directly importable. In this case, you should fully document every step you took to create the importable data in a subsection like this one.

### 3.1.3 Variable descriptions

Create a bullet list with the name of each variable in the dataset followed by any information the user would need to understand it.

- **variable\_\_name:** Variable description.
- **variable\_\_name2:** Description of second variable.

### 3.1.4 Data import code and summary

Once you've described the variables, enter an R chunk by selecting Code -> Insert Chunk, or Ctrl+Alt+I, give it a name to describe the dataset you are importing. After importing, export a dataframe summary using the command.

```
years = c("2003","2004")

download_and_read_data <- function(filename){

  if (!file.exists(file.path("raw-data",str_c(filename, ".csv")))) {

    destfile = file.path("raw-data",str_c(filename, ".zip"))

    download.file(url = str_c("https://www.tceq.texas.gov/assets/public/compliance/monops/air/ozonehist.",
                              str_c(filename, ".zip")), destfile, mode = "wb")

    unzip(destfile, exdir = "raw-data", junkpaths = T)
  }
  this_data <- read_csv(file.path("raw-data", str_c(filename, ".csv")))
}

oz_data <- lapply(str_c("oz_",years), download_and_read_data) %>% bind_rows()

## Parsed with column specification:
## cols(
##   .default = col_double(),
##   date = col_character(),
##   OZunit = col_character(),
##   OZmetxt = col_character()
## )

## See spec(...) for full column specifications.

## Warning: 1082 parsing failures.
##   row    col expected actual          file
## 17625 OZ1hrvh a double    ** 'raw-data/oz_2003.csv'
## 17626 OZ1hrvh a double    ** 'raw-data/oz_2003.csv'
## 17627 OZ1hrvh a double    ** 'raw-data/oz_2003.csv'
## 17628 OZ1hrvh a double    ** 'raw-data/oz_2003.csv'
## 17629 OZ1hrvh a double    ** 'raw-data/oz_2003.csv'
## .....
## See problems(...) for more details.
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   date = col_character(),
##   OZunit = col_character(),
##   OZmetxt = col_character(),
##   poc = col_logical(),
##   units = col_logical()
## )
## See spec(...) for full column specifications.

## Warning: 2609 parsing failures.
##   row    col                expected actual          file
## 17435 OZ1hrvh a double                ** 'raw-data/oz_2004.csv'
## 17435 units  1/0/T/F/TRUE/FALSE    ppb 'raw-data/oz_2004.csv'
## 17436 OZ1hrvh a double                ** 'raw-data/oz_2004.csv'
## 17436 units  1/0/T/F/TRUE/FALSE    ppb 'raw-data/oz_2004.csv'
## 17437 OZ1hrvh a double                ** 'raw-data/oz_2004.csv'
## .....
## See problems(...) for more details.
```

```
co_data <- lapply(str_c("co_",years), download_and_read_data) %>% bind_rows()
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   date = col_character(),
##   COunit = col_character(),
##   COmetxt = col_character()
## )
## See spec(...) for full column specifications.
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   date = col_character(),
##   COunit = col_character(),
##   COmetxt = col_character()
## )

## See spec(...) for full column specifications.
```

```
dfSummary(oz_data)
```

### 3.1.5 Data Frame Summary

```
oz_data
Dimensions: 51150 x 64
Duplicates: 0
```



No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
7	OZ1hr2 [numeric]	Mean (sd) : 18.6 (13) min < med < max: 0 < 18 < 81 IQR (CV) : 19 (0.7)	80 distinct values	: : : : : : : : : : : : : : : : : : :	878 (1.72%)
8	OZ1hr3 [numeric]	Mean (sd) : 17.7 (12.6) min < med < max: 0 < 17 < 81 IQR (CV) : 19 (0.7)	77 distinct values	: : : : : : : : : : : : . : : : : : .	882 (1.72%)
9	OZ1hr4 [numeric]	Mean (sd) : 16.4 (12.2) min < med < max: 0 < 15 < 75 IQR (CV) : 19 (0.7)	73 distinct values	: : : . : : : : : : : : . : : : : : .	864 (1.69%)
10	OZ1hr5 [numeric]	Mean (sd) : 14.5 (11.8) min < med < max: 0 < 13 < 74 IQR (CV) : 19 (0.8)	74 distinct values	: : : : : : : : : : : : : : .	867 (1.7%)
11	OZ1hr6 [numeric]	Mean (sd) : 13.5 (11.2) min < med < max: 0 < 11 < 74 IQR (CV) : 17 (0.8)	72 distinct values	: : : : . : : : . : : : : : .	889 (1.74%)
12	OZ1hr7 [numeric]	Mean (sd) : 15.8 (11.2) min < med < max: 0 < 14 < 81 IQR (CV) : 16 (0.7)	76 distinct values	\ . : . : : : : : : : : : .	937 (1.83%)
13	OZ1hr8 [numeric]	Mean (sd) : 21.4 (12.3) min < med < max: 0 < 20 < 96 IQR (CV) : 17 (0.6)	87 distinct values	: : : : : : . : : : : : : : : : .	1049 (2.05%)
14	OZ1hr9 [numeric]	Mean (sd) : 27.5 (13.9) min < med < max: 0 < 26 < 133 IQR (CV) : 18 (0.5)	104 distinct values	: : . : : : : : : : : : : .	1279 (2.5%)
15	OZ1hr10 [numeric]	Mean (sd) : 32.9 (15.6) min < med < max: 0 < 31 < 165 IQR (CV) : 20 (0.5)	127 distinct values	: : : : : : : . : : : : .	1348 (2.64%)
16	OZ1hr11 [numeric]	Mean (sd) : 37 (16.9) min < med < max: 0 < 34 < 196 IQR (CV) : 22 (0.5)	137 distinct values	: : : : : : : : : : : .	1378 (2.69%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
17	OZ1hr12 [numeric]	Mean (sd) : 39.9 (17.9) min < med < max: 0 < 37 < 195 IQR (CV) : 22 (0.4)	147 distinct values	: : . : : : : . : : : . .	1351 (2.64%)
18	OZ1hr13 [numeric]	Mean (sd) : 41.7 (18.5) min < med < max: 0 < 39 < 194 IQR (CV) : 23 (0.4)	156 distinct values	: . : : : : : : . : : : . .	1401 (2.74%)
19	OZ1hr14 [numeric]	Mean (sd) : 42.4 (18.8) min < med < max: 0 < 40 < 199 IQR (CV) : 24 (0.4)	158 distinct values	: . : : : : : : . : : : . .	1361 (2.66%)
20	OZ1hr15 [numeric]	Mean (sd) : 42.1 (18.8) min < med < max: 0 < 40 < 229 IQR (CV) : 24 (0.4)	156 distinct values	: : : : : : : : : : .	1158 (2.26%)
21	OZ1hr16 [numeric]	Mean (sd) : 40.4 (18.7) min < med < max: 0 < 38 < 176 IQR (CV) : 25 (0.5)	150 distinct values	: : : : : : . : : : : : : : : .	996 (1.95%)
22	OZ1hr17 [numeric]	Mean (sd) : 36.6 (18.7) min < med < max: 0 < 35 < 155 IQR (CV) : 25 (0.5)	141 distinct values	: . : : : : : : : : : : : : : .	921 (1.8%)
23	OZ1hr18 [numeric]	Mean (sd) : 31.1 (17.9) min < med < max: 0 < 30 < 131 IQR (CV) : 25 (0.6)	123 distinct values	: : : : . : : : : : : : : . : : : : : .	833 (1.63%)
24	OZ1hr19 [numeric]	Mean (sd) : 26.3 (16.3) min < med < max: 0 < 25 < 114 IQR (CV) : 23 (0.6)	108 distinct values	: . : : : : : : : : : : : . : : : : : .	827 (1.62%)
25	OZ1hr20 [numeric]	Mean (sd) : 23.7 (15.3) min < med < max: 0 < 23 < 119 IQR (CV) : 22 (0.6)	100 distinct values	: . . : : : : : : . : : : : : : : : : .	844 (1.65%)
26	OZ1hr21 [numeric]	Mean (sd) : 22 (14.7) min < med < max: 0 < 21 < 110 IQR (CV) : 21 (0.7)	94 distinct values	\ : . : : : : . : : : : : : : .	2191 (4.28%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
27	OZ1hr22 [numeric]	Mean (sd) : 21.2 (14.5) min < med < max: 0 < 20 < 91 IQR (CV) : 21 (0.7)	89 distinct values	\ . . : : . : : : : : : : : : : : :	7878 (15.4%)
28	OZ1hr23 [numeric]	Mean (sd) : 20.5 (14.1) min < med < max: 0 < 20 < 98 IQR (CV) : 21 (0.7)	83 distinct values	\ : : : : : : : : : : . : : : : .	2560 (5%)
29	OZunit [character]	1. ppb	51150 (100.0%)	IIIIIIIIIIIIIIIIIIII	0 (0%)
30	OZmetxt [character]	1. INSTRUMENTAL ULTRA VIOLET 2. INSTRUMENTAL ULTRA VIOLET	47899 (98.6%) 701 ( 1.4%)	IIIIIIIIIIIIIIIIIIII	2550 (4.99%)
31	OZmeth [numeric]	Mean (sd) : 55.4 (6.7) min < med < max: 19 < 56 < 87 IQR (CV) : 0 (0.1)	19 : 1066 ( 2.2%) 47 : 1037 ( 2.1%) 53 : 425 ( 0.9%) 56 : 45371 (93.4%) 87 : 701 ( 1.4%)	IIIIIIIIIIIIIIIIIIII	2550 (4.99%)
32	OZ1hrvh [numeric]	Mean (sd) : 23.6 (1.5) min < med < max: 1 < 24 < 99 IQR (CV) : 0 (0.1)	30 distinct values	: : : : :	2810 (5.49%)
33	OZ1hrvd [numeric]	Min : 0 Mean : 1 Max : 1	0 : 1066 ( 2.1%) 1 : 50084 (97.9%)	IIIIIIIIIIIIIIIIIIII	0 (0%)
34	OZ1hrpk [numeric]	Mean (sd) : 46.8 (19.4) min < med < max: 0 < 44 < 229 IQR (CV) : 24 (0.4)	179 distinct values	: : . : : : : : : : :	583 (1.14%)
35	OZ1hrav [numeric]	Mean (sd) : 26.8 (11.5) min < med < max: 0 < 25.5 < 91.6 IQR (CV) : 15.9 (0.4)	4402 distinct values	: : : : : : : : : : : : : : : : :	583 (1.14%)



No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
36	oz8hr0 [numeric]	Mean (sd) : 17 (11.4) min < med < max: 0 < 15.9 < 74.2 IQR (CV) : 16.6 (0.7)	641 distinct values	\ : : : : : : : : : : : : : : : : .	889 (1.74%)
37	oz8hr1 [numeric]	Mean (sd) : 17.1 (11) min < med < max: 0 < 15.9 < 75.9 IQR (CV) : 15.8 (0.6)	697 distinct values	. : . : : : : : : : : : : : . : : : : : .	894 (1.75%)
38	oz8hr2 [numeric]	Mean (sd) : 18.2 (10.6) min < med < max: 0 < 16.9 < 77.5 IQR (CV) : 15 (0.6)	767 distinct values	. : : : : : : : : : : : . : : : : : .	922 (1.8%)
39	oz8hr3 [numeric]	Mean (sd) : 19.9 (10.5) min < med < max: 0 < 18.8 < 79.1 IQR (CV) : 14.4 (0.5)	846 distinct values	. : : : . : : : : : : : . : : : : : .	962 (1.88%)
40	oz8hr4 [numeric]	Mean (sd) : 22.3 (10.6) min < med < max: 0 < 21.1 < 83.6 IQR (CV) : 14.7 (0.5)	896 distinct values	. : : : . : : : : : : : : : : : : .	1070 (2.09%)
41	oz8hr5 [numeric]	Mean (sd) : 25.3 (11.2) min < med < max: 0 < 24 < 86.8 IQR (CV) : 15.8 (0.4)	940 distinct values	. : : : : : : : . : : : : : : : : : .	1135 (2.22%)
42	oz8hr6 [numeric]	Mean (sd) : 28.7 (12.4) min < med < max: 0 < 27.1 < 93.5 IQR (CV) : 17.1 (0.4)	1063 distinct values	. : : : : : : . : : : : : : : : : .	1207 (2.36%)
43	oz8hr7 [numeric]	Mean (sd) : 32.3 (13.9) min < med < max: 0 < 30.4 < 105.6 IQR (CV) : 19 (0.4)	1195 distinct values	. : : . : : : . : : : : : : : : : .	1274 (2.49%)
44	oz8hr8 [numeric]	Mean (sd) : 35.6 (15.2) min < med < max: 0 < 33.4 < 122.5 IQR (CV) : 20.4 (0.4)	1317 distinct values	. : : . : : : : : : : : : : : : .	1308 (2.56%)
45	oz8hr9 [numeric]	Mean (sd) : 38 (16.3) min < med < max: 0 < 35.8 < 133.6 IQR (CV) : 21.6 (0.4)	1388 distinct values	. : : : : : : : : . : : : : : .	1327 (2.59%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
46	oz8hr10 [numeric]	Mean (sd) : 39.1 (17) min < med < max: 0 < 36.9 < 141.4 IQR (CV) : 22.6 (0.4)	1437 distinct values	: .: ::: :::: .....	1275 (2.49%)
47	oz8hr11 [numeric]	Mean (sd) : 38.9 (17.2) min < med < max: 0 < 36.5 < 141.3 IQR (CV) : 23.1 (0.4)	1432 distinct values	: .: ::: :::: .....	1204 (2.35%)
48	oz8hr12 [numeric]	Mean (sd) : 37.6 (16.9) min < med < max: 0 < 35.4 < 140.3 IQR (CV) : 22.9 (0.5)	1352 distinct values	: .: ::: :::: .....	1123 (2.2%)
49	oz8hr13 [numeric]	Mean (sd) : 35.5 (16.3) min < med < max: 0 < 33.5 < 132 IQR (CV) : 22.2 (0.5)	1257 distinct values	: .: ::: :::: .....	1033 (2.02%)
50	oz8hr14 [numeric]	Mean (sd) : 33.1 (15.7) min < med < max: 0 < 31.2 < 121.3 IQR (CV) : 21.6 (0.5)	1276 distinct values	: .: ::: :::: .....	968 (1.89%)
51	oz8hr15 [numeric]	Mean (sd) : 30.7 (15.1) min < med < max: 0 < 29 < 114 IQR (CV) : 21 (0.5)	1419 distinct values	: .: ::: :::: .....	924 (1.81%)
52	oz8hr16 [numeric]	Mean (sd) : 28 (14.6) min < med < max: 0 < 26.5 < 107.3 IQR (CV) : 20.4 (0.5)	1404 distinct values	: .: ::: :::: .....	873 (1.71%)
53	oz8hr17 [numeric]	Mean (sd) : 25.4 (14.1) min < med < max: 0 < 24.1 < 99.2 IQR (CV) : 19.8 (0.6)	1305 distinct values	: .: ::: :::: .....	1155 (2.26%)
54	oz8hr18 [numeric]	Mean (sd) : 23.1 (13.6) min < med < max: 0 < 22.1 < 92.6 IQR (CV) : 19.3 (0.6)	1232 distinct values	: .: ::: :::: .....	1212 (2.37%)
55	oz8hr19 [numeric]	Mean (sd) : 21.5 (13.2) min < med < max: 0 < 20.5 < 86.8 IQR (CV) : 18.9 (0.6)	1172 distinct values	: .: ::: :::: .....	1205 (2.36%)

No	Variable	Stats / Values	Freqs (% of Valid)	Graph	Missing
56	oz8hr20 [numeric]	Mean (sd) : 20.4 (12.8) min < med < max: 0 < 19.4 < 82.2 IQR (CV) : 18.6 (0.6)	1129 distinct values	. : : : : : : : : : . : : : : : : : : : : .	1206 (2.36%)
57	oz8hr21 [numeric]	Mean (sd) : 19.5 (12.5) min < med < max: 0 < 18.5 < 79.9 IQR (CV) : 18.2 (0.6)	1106 distinct values	. . : : : : : : : : : . : : : : : : : : : : .	1189 (2.32%)
58	oz8hr22 [numeric]	Mean (sd) : 18.5 (12.2) min < med < max: 0 < 17.5 < 77.9 IQR (CV) : 17.8 (0.7)	1056 distinct values	\ . : : : : : : : : : : : : : : : : .	1173 (2.29%)
59	oz8hr23 [numeric]	Mean (sd) : 17.6 (11.8) min < med < max: 0 < 16.6 < 75.8 IQR (CV) : 17.4 (0.7)	830 distinct values	\ . : : : . : : : : : : : : : : : : .	1166 (2.28%)
60	oz8hrpk [numeric]	Mean (sd) : 40.5 (16.7) min < med < max: 0 < 37.9 < 141.4 IQR (CV) : 22.1 (0.4)	1518 distinct values	: : : : : : : : . . : : : : .	580 (1.13%)
61	oz8hrvh [numeric]	Mean (sd) : 23.5 (3) min < med < max: 0 < 24 < 24 IQR (CV) : 0 (0.1)	25 distinct values	: : : : :	0 (0%)
62	oz8hrvd [numeric]	Min : 0 Mean : 1 Max : 1	0 : 1562 ( 3.0%) 1 : 49588 (97.0%)	IIIIIIIIIIIIIIIIIIII	0 (0%)
63	poc [logical]	1. TRUE	1464 (100.0%)	IIIIIIIIIIIIIIIIIIII	49686 (97.14%)
64	units [logical]	All NA's			51150 (100%)

```
export_summary_table(dfSummary(dataset_name))
```

While it will make your resulting file long, you should not modify the chunk options to suppress printing of code and output. I would likely not include this in the documentation for an actual paper I was submitting, but including them here will let me read your code and the output message from R and may help identify data import concerns early in the process. Since these files will exist only electronically, their length is less of a concern. If you like to print out files to proofread and want me to help you shorten the printed versions, let me know. We can temporarily modify the chunk options for printing and restore them before you submit the assignment.

## 4 Data Processing and Combination

*This section should include a discussion of the processing and merging steps needed to create your basic data. The code to implement these steps should be included in chunks in this section. Once the final merged data has been created, you should use the `dfSummary` function again to summarize the data you will be using. You should also save a file containing all the objects you will use in your final analysis to the `processed_data` folder.*

## 5 Analysis Variables

This section should include a description of all the variables that are used in your final analysis. At the end of the section, you should save all of these variables in the `processed_data` folder of your repository.

## 6 Discussion of Data

*This section should include a discussion of any data patterns you notice based on the summaries created in the code above.*