

# Sample Datasets Info Page

Sean Conway

8/17/2021

## Datasets

This document summarizes the datasets that have been collected for use in DACSS 601 for the August 2021 session. All files can be found in the `_data` folder on the course blog. Note that some of these datasets require significant wrangling/cleaning. Also note that any `.xls/.xlsx` files may have multiple sheets, so it will be helpful to open these files in a spreadsheet software first, to examine the file you are reading in.

Also note that there are both **basic** and **advanced** versions of the datasets. The basic versions are clean and can be read into R fairly easily. These are great to practice on if you are new to R or need a refresher on importing data. For those who are more advanced R programmers, the advanced datasets will require significant work to be imported into R and tidied.

## Hotel Bookings

This dataset contains hotel bookings from 2015-2017. Each row is an individual hotel booking. This dataset is **only available as an advanced dataset**. The file is named `hotel_bookings.csv`. Because the file format is `.csv`, we can use the function `read_csv()` from the `readr` package to read in the data to R.

```
hotels <- read_csv(here("R", "data", "hotel_bookings.csv"))
hotels
```

```
## # A tibble: 119,390 x 32
##   hotel          is_canceled lead_time arrival_date_year arrival_date_month
##   <chr>          <dbl>      <dbl>         <dbl> <chr>
## 1 Resort Hotel      0        342          2015 July
## 2 Resort Hotel      0        737          2015 July
## 3 Resort Hotel      0         7          2015 July
## 4 Resort Hotel      0        13          2015 July
## 5 Resort Hotel      0        14          2015 July
## 6 Resort Hotel      0        14          2015 July
## 7 Resort Hotel      0         0          2015 July
## 8 Resort Hotel      0         9          2015 July
## 9 Resort Hotel      1        85          2015 July
## 10 Resort Hotel     1        75          2015 July
## # ... with 119,380 more rows, and 27 more variables:
## #   arrival_date_week_number <dbl>, arrival_date_day_of_month <dbl>,
## #   stays_in_weekend_nights <dbl>, stays_in_week_nights <dbl>, adults <dbl>,
## #   children <dbl>, babies <dbl>, meal <chr>, country <chr>,
## #   market_segment <chr>, distribution_channel <chr>, is_repeated_guest <dbl>,
## #   previous_cancellations <dbl>, previous_bookings_not_canceled <dbl>,
## #   reserved_room_type <chr>, assigned_room_type <chr>, ...
```

Source: <https://www.kaggle.com/jessemostipak/hotel-booking-demand>

Also see the link for a detailed key.

## 2019 New York City Air BnB Bookings

This dataset contains Air Bnb bookings from 2019 in New York City. Each row contains an individual Air Bnb listing, and each column contains information about it (e.g., number of reviews per month, price, data of last review). This dataset is **only available as an advanced dataset**. The file is named `AB_NYC_2019.csv`. Because the file format is `.csv`, we can use the function `read_csv()` from the `readr` package to read in the data to R.

```
air_bnb <- read_csv(here("R", "data", "AB_NYC_2019.csv"))
air_bnb
```

```
## # A tibble: 48,895 x 16
##       id name          host_id host_name neighbourhood_g~ neighbourhood latitude
##   <dbl> <chr>          <dbl> <chr>      <chr>          <chr>          <dbl>
## 1  2539 Clean & qui~    2787 John      Brooklyn    Kensington    40.6
## 2  2595 Skylit Midt~    2845 Jennifer  Manhattan   Midtown        40.8
## 3  3647 THE VILLAGE~    4632 Elisabeth  Manhattan   Harlem          40.8
## 4  3831 Cozy Entire~    4869 LisaRoxan~ Brooklyn    Clinton Hill    40.7
## 5  5022 Entire Apt:~    7192 Laura      Manhattan   East Harlem    40.8
## 6  5099 Large Cozy ~    7322 Chris       Manhattan   Murray Hill    40.7
## 7  5121 BlissArtsSp~    7356 Garon       Brooklyn    Bedford-Stuy~    40.7
## 8  5178 Large Furni~    8967 Shunichi    Manhattan   Hell's Kitch~    40.8
## 9  5203 Cozy Clean ~    7490 MaryEllen  Manhattan   Upper West S~    40.8
## 10 5238 Cute & Cozy~    7549 Ben        Manhattan   Chinatown       40.7
## # ... with 48,885 more rows, and 9 more variables: longitude <dbl>,
## #   room_type <chr>, price <dbl>, minimum_nights <dbl>,
## #   number_of_reviews <dbl>, last_review <date>, reviews_per_month <dbl>,
## #   calculated_host_listings_count <dbl>, availability_365 <dbl>
```

```
glimpse(air_bnb)
```

```
## Rows: 48,895
## Columns: 16
## $ id          <dbl> 2539, 2595, 3647, 3831, 5022, 5099, 512~
## $ name        <chr> "Clean & quiet apt home by the park", "~
## $ host_id     <dbl> 2787, 2845, 4632, 4869, 7192, 7322, 735~
## $ host_name   <chr> "John", "Jennifer", "Elisabeth", "LisaR~
## $ neighbourhood_group <chr> "Brooklyn", "Manhattan", "Manhattan", "~
## $ neighbourhood <chr> "Kensington", "Midtown", "Harlem", "Cli~
## $ latitude    <dbl> 40.64749, 40.75362, 40.80902, 40.68514,~
## $ longitude   <dbl> -73.97237, -73.98377, -73.94190, -73.95~
## $ room_type   <chr> "Private room", "Entire home/apt", "Pri~
## $ price       <dbl> 149, 225, 150, 89, 80, 200, 60, 79, 79,~
## $ minimum_nights <dbl> 1, 1, 3, 1, 10, 3, 45, 2, 2, 1, 5, 2, 4~
## $ number_of_reviews <dbl> 9, 45, 0, 270, 9, 74, 49, 430, 118, 160~
## $ last_review  <date> 2018-10-19, 2019-05-21, NA, 2019-07-05~
## $ reviews_per_month <dbl> 0.21, 0.38, NA, 4.64, 0.10, 0.59, 0.40,~
## $ calculated_host_listings_count <dbl> 6, 2, 1, 1, 1, 1, 1, 1, 1, 4, 1, 1, 3, ~
## $ availability_365 <dbl> 365, 355, 365, 194, 0, 129, 0, 220, 0, ~
```

Source: <https://www.kaggle.com/dgomonov/new-york-city-airbnb-open-data>

Also see the link for a detailed key.

## 2017 Australian Marriage Law

Data on public opinion of a proposed same sex marriage law in Australia in 2017. The **basic** version of this dataset can be found as `australian_marriage_tidy.csv` and `australian_marriage_tidy.xlsx`. The advanced version is `australian_marriage_law_postal_survey_2017_-_response_final.xlsx`, so we can use the function `read_excel()` to read in the data. However, this advanced dataset was designed as an Excel spreadsheet, and so will take some extra work to be read into R.

Source: <https://www.abs.gov.au/ausstats/abs@.nsf/mf/1800.0>

## DOD Active Duty Marital Status

Count data on various demographic characteristics, notably marital status and child status, by pay grade, for multiple branches of the military (as well as DOD as a whole). This dataset is **only available as an advanced dataset**. This file is called `ActiveDuty_MaritalStatus.xls`. However, this dataset was designed as an Excel spreadsheet, and so will take some extra work to be read into R.

Source: <https://catalog.data.gov/dataset/active-duty-marital-status/resource/638cad03-b16c-48ac-8346-f858ff89d202>

## Public School Characteristics 2017-2018

Data on characteristics of every US public school from 2017-2018. File is called `Public_School_Characteristics_2017-18.csv`. Note that this file is fairly large, and if you aren't careful, you may encounter parsing errors when reading in the file.

Source: <https://catalog.data.gov/dataset/public-school-characteristics-2017-18>

## 2012 US Railroad Employment.

Data breaking down US railroad employment numbers in 2012 by state and county. The **basic** versions are divided into county data and state data. The **basic** files are `railroad_2012_clean_county.csv` and `railroad_2012_clean_county.xlsx`, and `railroad_2012_clean_state.csv` and `railroad_2012_clean_state.xlsx`. The **advanced** file is `StateCounty2012.xls`.

Source: <https://catalog.data.gov/dataset/total-railroad-employment-by-state-and-county-2012/resource/5a0b2831-23b9-4ce9-82e9-87a7d8f2c5d8>

## Organic Egg & Poultry Prices

Data on organic egg & poultry prices in the US from 2004-2013. The basic versions of the files are `poultry_tidy.csv` and `poultry_tidy.xlsx`, as well as `eggs_tidy.csv` and `eggs_tidy.xlsx`. The **advanced** file is `organiceggpoultry.xls`.

Source: <https://www.ers.usda.gov/data-products/organic-prices.aspx>