Data types

- 1) Provide a URL to the dataset. I downloaded my dataset from http://data.insideairbnb.com/united-states/ca/san-francisco/2017-10-02/data/listings.csv.gz (http://data.insideairbnb.com/united-states/ca/san-francisco/2017-10-02/data/listings.csv.gz)
- 2) Explain why you chose this dataset.

I have used this dataset before for a project i was interested in.

- 3) What are the entities in this dataset? How many are there? The entities are specific Airbnb listings in San Francisco. There are 8706 listings.
- 4) How many attributes are there in this dataset? There are 95 attributes.
- 5) What is the datatype of each attribute (categorical -ordered or unordered-, numeric -discrete or continuous-, datetime, geolocation, other)? Write a short sentence stating how you determined the type of each attribute. Do this for at least 5 attributes, if your dataset contains more than 10 attributes, choose 10 of them to describe.

Num	Name	Туре	Description
1	'ID'	Categorical	Identifier of each listing which is a finite set.
2	'Neighborhood'	Categorical Unordered	Finite set of number of neighborhoods
3	'Latitude'	Numerical Continuous	Infinite precision can be applied to yield an infinite set.
4	'Longitude'	Numerical Continuous	Infinite precision can be applied to yield an infinte set.
5	'bedrooms'	Numerical Discrete	Finite set of number of beds.
6	'Maximum nights'	Numerical Discrete	Finite set of nights listing allows.

1 of 5 2/6/2018, 10:13 PM

Num	Name	Туре	Description
7	'Host listings count'	Numerical Discrete	Hosts can only list a finite number of listings.
8	'Accommodates'	Numerical Discrete	listings can only accomoodate a finite number of people
9	'Bed Type'	Categorical Unordered	Finite set of bed types
10	'Number of Reviews'	Numerical Discrete	Finite set of reviews.

6) Write R code that loads the dataset using function <code>read_csv</code>. Were you able to load the data successfully? If no, why not? Yes but ther were some parsing errors.

```
library(tidyverse)
# loading code goes here
mydata <- read_csv("listings.csv")</pre>
```

```
mydata
```

2 of 5

```
# A tibble: 8,706 x 95
           id listing url scrape id last scraped name summary
##
                                                                 space
                             <dbl> <date> <chr> <chr>
       <int> <chr>
##
                                                                 <chr>
   1 1.10e7 https://ww~ 2.02e13 2017-04-02 Simp~ I signed ~ This is cl~
   2 8.05e6 https://ww~ 2.02e13 2017-04-02 Spac~ Enjoy the~ "We love o~
##
   3 1.40e7 https://ww~ 2.02e13 2017-04-02 Priv~ Awesome 1~ <NA>
##
##
   4 1.61e7 https://ww~ 2.02e13 2017-04-02
                                                Spac~ Quiet 1 b~ <NA>
   5 9.08e6 https://ww~ 2.02e13 2017-04-02
                                                Edwa~ Bedroom w~ <NA>
##
   6 9.36e5 https://ww~ 2.02e13 2017-04-02
                                                Sunn~ Come stay~ Located ju~
##
##
   7 7.05e6 https://ww~ 2.02e13 2017-04-02
                                                Bout~ Boutique ~ <NA>
   8 8.00e6 https://ww~ 2.02e13 2017-04-02
                                                Pret~ Quiet one~ <NA>
   9 1.28e7 https://ww~ 2.02e13 2017-04-02
                                                Spac~ Built in ~ <NA>
##
## 10 1.04e7 https://ww~ 2.02e13 2017-04-02 Brig~ Located i~ Located in~
## # ... with 8,696 more rows, and 88 more variables: description <chr>,
      experiences offered <chr>, neighborhood overview <chr>, notes <chr>,
## #
      transit <chr>, access <chr>, interaction <chr>, house rules <chr>,
## #
      thumbnail url <chr>, medium url <chr>, picture url <chr>,
      xl picture url <chr>, host id <int>, host url <chr>, host name <chr>,
####
## #
      host_since <date>, host_location <chr>, host_about <chr>,
      host response time <chr>, host response rate <chr>,
## #
## #
      host acceptance rate <chr>, host is superhost <chr>,
## #
      host thumbnail url <chr>, host picture url <chr>,
## #
      host_neighbourhood <chr>, host_listings_count <int>,
## #
      host total listings count <int>, host verifications <chr>,
      host has profile pic <chr>, host identity verified <chr>,
## #
       street <chr>, neighbourhood <chr>, neighbourhood cleansed <chr>,
## #
## #
       neighbourhood group cleansed <chr>, city <chr>, state <chr>,
       zipcode <int>, market <chr>, smart location <chr>, country code <chr>,
## #
       country <chr>, latitude <dbl>, longitude <dbl>,
## #
       is_location_exact <chr>, property_type <chr>, room_type <chr>,
## #
## #
       accommodates <int>, bathrooms <dbl>, bedrooms <int>, beds <int>,
## #
      bed type <chr>, amenities <chr>, square feet <int>, price <chr>,
## #
       weekly price <chr>, monthly price <chr>, security deposit <chr>,
## #
       cleaning fee <chr>, guests included <int>, extra people <chr>,
## #
       minimum nights <int>, maximum nights <int>, calendar updated <chr>,
## #
      has availability <chr>, availability 30 <int>, availability 60 <int>,
## #
      availability 90 <int>, availability 365 <int>,
## #
       calendar last scraped <date>, number of reviews <int>,
## #
       first review <date>, last review <date>, review scores rating <int>,
       review scores accuracy <int>, review scores cleanliness <int>,
## #
       review_scores_checkin <int>, review_scores communication <int>,
## #
## #
       review_scores_location <int>, review_scores_value <int>,
## #
       requires license <chr>, license <chr>, jurisdiction names <chr>,
## #
       instant bookable <chr>, cancellation policy <chr>,
####
       require guest profile picture <chr>,
## #
       require guest phone verification <chr>,
## #
       calculated host listings count <int>, reviews per month <dbl>
```

Wrangling

3 of 5

1)I am trying to find a listing on Airbnb that can accommodate at least a family of 4 (if it accommodates more, great I can bring more people but 4 is the minimum) that also has really good ratings. Price is not a concern here, so I left it off but ideally someone else would want to compare prices to get the best deal.

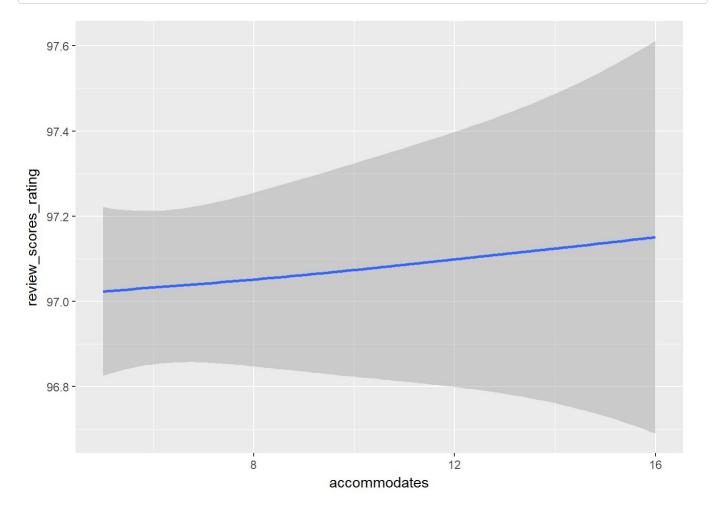
```
# pipeline goes here
data <- mydata %>%
select(id,review_scores_rating, accommodates) %>%
filter(review_scores_rating > 90, accommodates > 4) %>%
arrange(desc(review_scores_rating))
```

Plotting

1. I wanted to see if a review score rating was based on how many people the listing accommodated and from my results it does not appear so for review scores above 90 and accommodations above 4. On the X-axis I have the different accommodations of listings and I have the dependent variable on the Y-Axis (review score)

```
# plot goes here
data %>%
ggplot(mapping=aes(x = accommodates, y = review_scores_rating)) + geom_smooth()
```

```
## `geom_smooth()` using method = 'gam'
```



4 of 5 2/6/2018, 10:13 PM

5 of 5