

# DATA201 Group Project

Taking an inside look at AirBNBs around the world!

---

By Yazeed, Heran and Chaarvee



# Our research question and what we wanted to achieve.

---

- Comparing AirBNB prices around the world and if the countries population and GDP per capita has any relation to the AirBNB prices?
- We also wanted to compare the number of bedrooms and bathrooms for each AirBNB and to see if that also has any relation to the price of AirBNBs.



# Why we chose AirBNBS?

---

- As a group we were all interested in this idea and felt like it was something unique, and maybe it could help us choose our next travel destinations.

# How we acquired the data for AirBNBs?

- We downloaded CSVs for each of the countries we wanted to scrape from the following website <http://insideairbnb.com/get-the-data/>
- We then read in the CSVs for the selected countries
- Here is an example of the raw data that the CSVs spat out

```
Rows: 46,412
Columns: 80
$ id
$ listing_url
$ scrape_id
$ last_searched
$ last_scraped
$ name
$ description
$ neighborhood_overview
$ picture_url
$ host_id
$ host_url
$ host_name
$ host_since
$ host_location
$ host_about
$ host_response_time
$ host_response_rate
$ host_acceptance_rate
$ host_is_superhost
$ host_thumbnail_url
$ host_picture_url
$ host_neighbourhood
$ host_listings_count
$ host_total_listings_count
$ host_verifications
$ host_has_profile_pic
$ host_identity_verified
$ neighbourhood
$ latitude
$ longitude
$ property_type
$ room_type
$ accommodates
$ bathrooms
$ bathrooms_text
$ bedrooms
$ beds
$ amenities
$ price
$ minimum_nights
$ maximum_nights
$ minimum_minimum_nights
$ maximum_minimum_nights
$ minimum_maximum_nights
$ maximum_maximum_nights
$ minimum_nights_avg_ntm
$ maximum_nights_avg_ntm
$ calendar_updated
$ has_availability

<dbl> 6113, 35325, 46071, 48443...
<chr> "https://www.airbnb.com/r...
<dbl> 2.02309e+13, 2.02309e+13,...
<date> NA, NA, 2023-09-03, NA, ...
<date> 2023-09-04, 2023-09-04, ...
<chr> "Place to stay in Otaki '...
<chr> "<b>The space</b><br />La...
<chr> NA, "We are located in th...
<chr> "https://a0.muscache.com/...
<dbl> 12177, 152089, 202747, 22...
<chr> "https://www.airbnb.com/u...
<chr> "Dianne", "Chika", "Donna...
<date> 2009-04-03, 2010-06-25, ...
<chr> "Otaki, New Zealand", "Qu...
<chr> "I have a great interest ...
<chr> "N/A", "within a few hour...
<chr> "N/A", "100%", "100%", "N...
<chr> "N/A", "94%", "100%", "N/...
<lgl> FALSE, TRUE, FALSE, FALSE...
<chr> "https://a0.muscache.com/...
<chr> "https://a0.muscache.com/...
<chr> NA, NA, NA, NA, NA, NA, N...
<dbl> 1, 2, 1, 2, 2, 1, 1, 1, 1...
<dbl> 1, 3, 2, 3, 3, 1, 3, 1, 1...
<chr> "[email]", "phone"]", "[...
<lgl> TRUE, TRUE, TRUE, FALSE, ...
<lgl> FALSE, TRUE, TRUE, FALSE,...
<chr> NA, "Queenstown, Otago, N...
<dbl> -40.75807, -45.00532, -38...
<dbl> 175.1564, 168.7771, 175.7...
<chr> "Private room", "Entire h...
<chr> "Private room", "Entire h...
<dbl> 2, 6, 8, 1, 2, 4, 1, 4, 2...
<lgl> NA, NA, NA, NA, NA, NA, N...
<chr> NA, "1.5 baths", "2 baths...
<dbl> 1, 3, 5, 1, 1, 1, 1, NA, ...
<dbl> 1, 6, 7, NA, 1, 2, NA, 2,...
<chr> "[Breakfast]", "Wifi"...
<chr> "$109.00", "$250.00", "$2...
<dbl> 1, 13, 2, 1, 1, 2, 1, 1, ...
<dbl> 21, 149, 10, 730, 730, 73...
<dbl> 1, 13, 2, 1, 1, 2, 1, 1, ...
<dbl> 1, 13, 2, 1, 1, 2, 1, 1, ...
<dbl> 21, 149, 10, 730, 730, 73...
<dbl> 21, 149, 10, 730, 730, 73...
<dbl> 1.0, 13.0, 2.0, 1.0, 1.0,...
<dbl> 21, 149, 10, 730, 730, 73...
<lgl> NA, NA, NA, NA, NA, NA, N...
<lgl> FALSE, TRUE, TRUE, FALSE,...
```



# Wrangling the data

- The CSVs for each AirBNB location had thousands of variables we could select from.
- For this project we only wrangled the price, number of bathrooms and the number of bedrooms
- Here is a small example of the wrangled data

A tibble: 87946 × 3

price	bathrooms_text	bedrooms
<chr>	<chr>	<dbl>
\$42.00	1.5 shared baths	NA
\$175.00	1 bath	2
\$79.00	1 shared bath	NA
\$150.00	1 bath	1
\$46.00	1 shared bath	NA
\$476.00	2 baths	3
\$371.00	2 baths	2
\$250.00	1.5 baths	1
\$75.00	1 bath	1
\$29.00	1.5 shared baths	NA
\$120.00	1 bath	2
\$50.00	1 shared bath	NA
\$151.00	1 bath	1
\$35.00	1 shared bath	NA
\$140.00	1 private bath	NA
\$90.00	1 bath	1



# How we acquired the data for world population and GDP per capita?

- We scraped this website <https://www.worldometers.info/gdp/gdp-by-country/>
- We selected the data we wanted from the website which was the Countries name, population and GDP per capita.

```
[32]: countries = world_data %>%  
      htl_nodes("example2.g") %>%  
      htl_text()  
      countries  
  
'United States', 'China', 'Japan', 'Germany', 'India', 'United Kingdom', 'France', 'Russia', 'Canada', 'Italy', 'Brazil', 'Australia', 'South Korea', 'Mexico', 'Spain', 'Indonesia', 'Saudi Arabia', 'Netherlands', 'Turkey',  
'Switzerland', 'Poland', 'Argentina', 'Sweden', 'Norway', 'Belgium', 'Ireland', 'Israel', 'United Arab Emirates', 'Thailand', 'Nigeria', 'Egypt', 'Austria', 'Singapore', 'Bangladesh', 'Vietnam', 'Malaysia', 'South Africa',  
'Philippines', 'Denmark', 'Iran', 'Pakistan', 'Hong Kong', 'Colombia', 'Romania', 'Chile', 'Czech Republic (Czechia)', 'Finland', 'Portugal', 'New Zealand', 'Algeria', 'Kuwait',  
'Hungary', 'Ukraine', 'Morocco', 'Ethiopia', 'Slovakia', 'Ecuador', 'Oman', 'Dominican Republic', 'Kenya', 'Angola', 'Guatemala', 'Luxembourg', 'Uzbekistan', 'Azerbaijan', 'Panama', 'Tanzania', 'Sri Lanka',  
'Ghana', 'Belarus', 'Uruguay', 'Croatia', 'Lithuania', 'Côte d'Ivoire', 'Costa Rica', 'Serbia', 'Slovenia', 'Myanmar', 'DR Congo', 'Sudan', 'Jordan', 'Tunisia', 'Libya', 'Uganda', 'Bahrain', 'Cameroon', 'Bolivia', 'Paraguay',  
'Latvia', 'Nepal', 'Estonia', 'El Salvador', 'Honduras', 'Papua New Guinea', 'Cambodia', 'Zambia', 'Cyprus', 'Trinidad and Tobago', 'Iceland', 'Senegal', 'Georgia', 'Bosnia and Herzegovina', 'Macao', 'Guinea', 'Gabon',  
'Zimbabwe', 'Botswana', 'Haiti', 'Armenia', 'State of Palestine', 'Burkina Faso', 'Albania', 'Mali', 'Mozambique', 'Malta', 'Benin', 'Jamaica', 'Mongolia', 'Brunei', 'Laos', 'Nicaragua', 'Guyana', 'Madagascar', 'Congo',  
'Moldova', 'Niger', 'North Macedonia', 'Bahrain', 'Malawi', 'Mauritius', 'Bahamas', 'Chad', 'Namibia', 'Equatorial Guinea', 'Kyrgyzstan', 'Tajikistan', 'Mauritania', 'Togo', 'Maldives', 'Montenegro', 'Barbados', 'Yemen',  
'Eswatini', 'Liberia', 'Sierra Leone', 'Suriname', 'Andorra', 'Timor-Leste', 'Burundi', 'Belize', 'Lesotho', 'Central African Republic', 'Cabo Verde', 'Gambia', 'Saint Lucia', 'Antigua and Barbuda', 'Guinea-Bissau',  
'Solomon Islands', 'Seychelles', 'Grenada', 'Comoros', 'Vanuatu', 'Saint Kitts & Nevis', 'St. Vincent & Grenadines', 'Samoa', 'Dominica', 'Sao Tome & Principe', 'Micronesia', 'Marshall Islands', 'Kiribati', 'Tuvalu',  
  
[33]: population = world_data %>%  
      htl_nodes("td:nth-child(6)") %>%  
      htl_text()  
      population  
  
'338,289,857', '1,426,887,337', '123,951,692', '83,368,843', '1,417,173,173', '87,508,936', '64,626,628', '144,713,314', '38,464,327', '59,037,474', '215,313,498', '26,177,411', '61,816,810', '127,504,125', '47,558,630',  
'278,501,339', '36,408,830', '17,864,014', '34,341,241', '36,740,472', '19,847,345', '45,310,318', '10,548,347', '4,434,319', '11,665,930', '5,023,109', '9,038,309', '9,441,129', '71,697,035', '218,541,212', '110,990,103',  
'8,939,617', '5,975,689', '171,186,372', '98,186,856', '33,938,221', '59,893,885', '115,559,009', '5,882,261', '88,550,570', '235,824,862', '7,488,868', '618,740,024', '19,659,267', '19,603,733', '10,493,986', '5,540,749',  
'44,496,122', '10,270,865', '5,185,288', '34,048,588', '2,695,122', '19,397,998', '10,384,977', '44,903,225', '4,268,873', '9,967,308', '39,701,739', '37,457,971', '123,379,924', '5,643,453', '18,001,000', '4,576,298',  
'11,228,821', '54,027,487', '35,588,987', '17,843,908', '6,781,953', '647,599', '34,627,652', '10,358,074', '4,408,581', '65,497,748', '21,832,143', '33,475,871', '9,534,954', '3,422,794', '4,030,358', '2,750,055',  
'28,160,542', '5,180,829', '2,721,365', '2,119,844', '54,179,306', '99,010,212', '46,874,204', '11,285,869', '12,356,117', '6,812,341', '47,249,585', '1,472,233', '27,914,536', '12,224,110', '6,780,744', '1,850,651',  
'30,547,580', '1,226,062', '6,336,392', '10,432,860', '10,142,819', '16,767,842', '20,017,670', '1,251,488', '17,318,449', '3,744,385', '3,233,528', '695,168', '13,858,341', '2,388,992', '16,320,537',  
'6,530,295', '11,584,096', '2,760,469', '5,250,072', '22,673,702', '2,842,321', '22,593,590', '32,968,518', '533,286', '13,352,854', '2,827,377', '3,268,366', '448,002', '7,529,705', '6,948,392', '908,735', '29,611,714',  
'5,970,424', '3,272,996', '26,207,877', '2,093,599', '13,776,698', '20,405,317', '1,999,469', '409,984', '17,723,315', '2,567,012', '1,674,908', '6,630,823', '9,952,787', '4,736,139', '8,848,699', '523,787', '527,082',  
'281,635', '929,766', '1,201,670', '5,302,681', '8,605,718', '618,040', '79,824', '1,341,296', '12,889,576', '405,272', '2,305,825', '5,579,144', '593,149', '2,705,992', '179,851', '93,763', '2,105,566', '724,273', '107,118',  
'125,438', '836,774', '326,740', '47,657', '103,848', '222,382', '72,737', '227,380', '539,013', '41,569', '131,232', '11,312',  
  
[35]: gdp_per_capita = world_data %>%  
      htl_nodes("td:nth-child(7)") %>%  
      htl_text()  
      gdp_per_capita  
  
'$75,269', '$12,598', '$34,135', '$48,845', '$2,389', '$45,485', '$43,061', '$15,482', '$55,646', '$34,053', '$8,918', '$64,003', '$32,138', '$11,091', '$20,385', '$4,788', '$30,438', '$56,429', '$10,616', '$92,410',  
'$17,268', '$13,904', '$55,543', '$106,594', '$49,640', '$105,362', '$57,758', '$53,758', '$6,909', '$2,164', '$4,295', '$52,732', '$2,688', '$4,164', '$1,972', '$6,776', '$3,699', '$67,220', '$4,388', '$1,597',  
'$48,050', '$6,630', '$15,324', '$15,355', '$27,723', '$50,684', '$5,937', '$24,530', '$47,680', '$7,126', '$88,046', '$11,373', '$21,095', '$4,274', '$43,233', '$17,938', '$4,043', '$3,582', '$1,028', '$20,461', '$6,391',  
'$25,057', '$10,121', '$2,099', '$2,999', '$5,324', '$13,129', '$127,046', '$2,322', '$7,600', '$17,358', '$1,156', '$3,408', '$2,176', '$7,634', '$20,795', '$17,608', '$25,576', '$2,486', '$13,199', '$8,794', '$29,303',  
'$1,096', '$686', '$1,102', '$4,205', '$3,777', '$6,716', '$964', '$30,152', '$1,588', '$3,523', '$6,153', '$1,337', '$28,732', '$5,127', '$3,040', '$3,020', '$1,787', '$1,488', '$22,724', '$18,222', '$74,663',  
'$1,599', '$6,071', '$7,685', '$31,618', '$1,532', '$8,820', '$1,267', '$2,738', '$1,748', '$7,014', '$3,640', '$833', '$6,643', '$833', '$541', '$33,319', '$1,303', '$6,047', '$4,947', '$3,192', '$2,089', '$2,259', '$18,960'
```



# Wrangling the data

- We then created a function which scrapes the entire webpage and only selects the countries we wanted their population and GDP per capita.

```
[44]: scrape_and_filter <- function(url) {  
  # Scrape data  
  webpage <- read_html(url)  
  
  countries <- webpage %>%  
    html_nodes('#example2 a') %>%  
    html_text()  
  
  population <- webpage %>%  
    html_nodes('td:nth-child(6)') %>%  
    html_text() %>%  
    gsub("[^0-9.]","", ".") %>%  
    as.numeric()  
  
  gdp_per_capita <- webpage %>%  
    html_nodes('td:nth-child(7)') %>%  
    html_text() %>%  
    gsub("[^0-9.]","", ".") %>%  
    as.numeric()  
  
  # Convert to data frame  
  data <- data.frame(  
    Country = countries,  
    Population = population,  
    GDP = gdp,  
    GDP_per_Capita = gdp_per_capita  
  )  
  
  # Filter for the required countries  
  filtered_data <- data %>%  
    filter(Country %in% c('United Kingdom', 'Japan', 'New Zealand', 'Greece', 'Belgium', 'Netherlands', 'United States', 'Spain', 'Thailand', 'Australia'))  
  
  return(filtered_data)  
}  
  
# Call the function  
url <- "https://www.worldometers.info/gdp/gdp-by-country/"  
result <- scrape_and_filter(url)  
print(result)
```

	Country	Population	GDP	GDP_per_Capita
1	United States	338289857	\$25,462,700,000,000	75269
2	Japan	123951692	\$4,231,140,000,000	34135
3	United Kingdom	67508936	\$3,070,670,000,000	45485
4	Australia	26177413	\$1,675,420,000,000	64083
5	Spain	47559638	\$1,397,510,000,000	29385
6	Netherlands	17564014	\$991,115,000,000	56429
7	Belgium	11655930	\$578,604,000,000	49640
8	Thailand	71697030	\$495,341,000,000	6909
9	New Zealand	5185288	\$247,234,000,000	47680
10	Greece	10384971	\$219,066,000,000	21095

# Difficulties we had during the project

---

- Something that was out of our control was some of the inputs for the number of bathrooms and bedrooms showed up as NA. We assume this is because the people who have listed up their homes on AirBNB have just left these sections blank for some odd reason.
- The pricing of AirBNBs being in the local currencies of the countries we had scraped (for example Japanese AirBNB data was in JPY and not NZD) which would make it much more difficult for us to answer our research question - however we did have a solution for this
- Lots and lots of raw data!



# How we overcame the issue pricing being in different currencies?

- After doing lots of research we found this website <https://fixer.io/> which has an API with the live exchange rates and helped us exchange the pricing to \$NZD which made it much easier for us to answer the research question.
- Here is the code for how we did this

```
[42]: # API key for fixer.io a live exchange rate site so we can convert the price column from JPY to NZD
      api_key <- "58c8eb71a98bc0e70380918c4ec1041c"

[43]: response <- GET(paste0("http://data.fixer.io/api/latest?access_key=", api_key, "&symbols=JPY,NZD"))

[44]: data <- fromJSON(content(response, as = "text"))

[45]: exchange_rate_jpy_to_nzd <- data$rates$NZD / data$rates$JPY

[46]: japan_airbnb$price_nzd <- japan_airbnb$price * exchange_rate_jpy_to_nzd

[47]: japan_airbnb_nzd <- japan_airbnb %>%
      mutate(price_nzd = price * exchange_rate_jpy_to_nzd) %>%
      select(price_nzd, bathrooms_text, bedrooms)

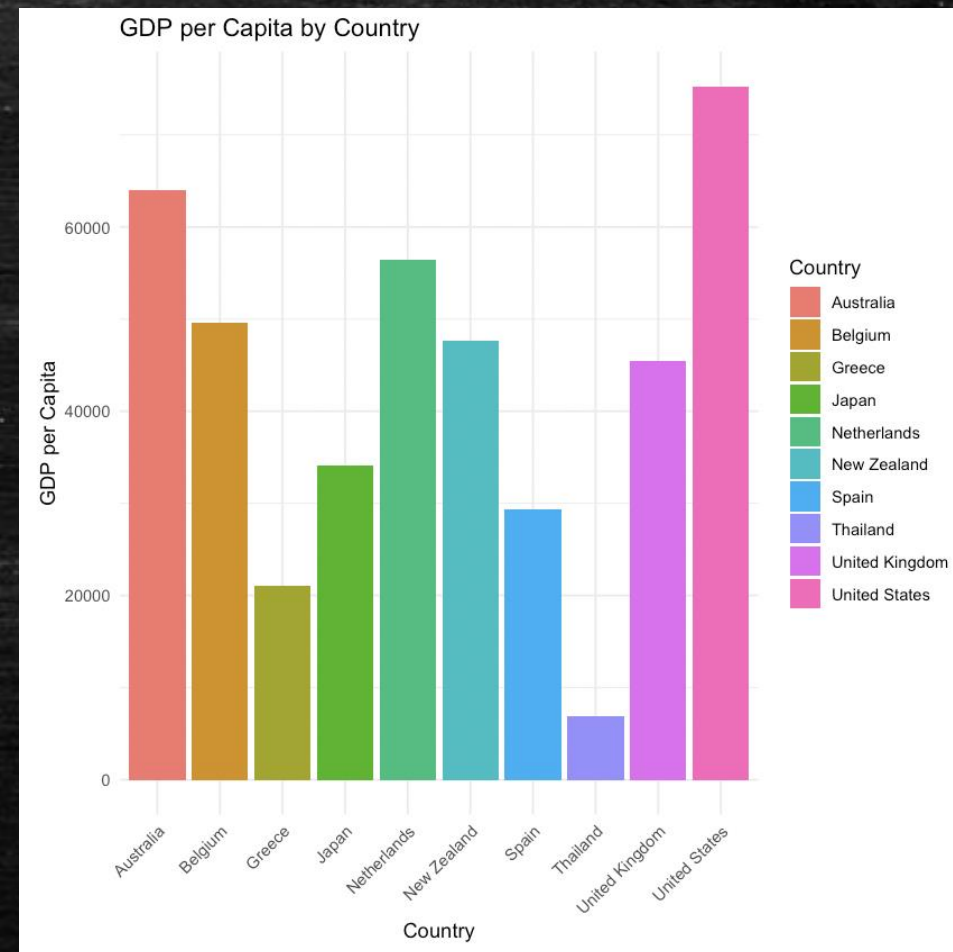
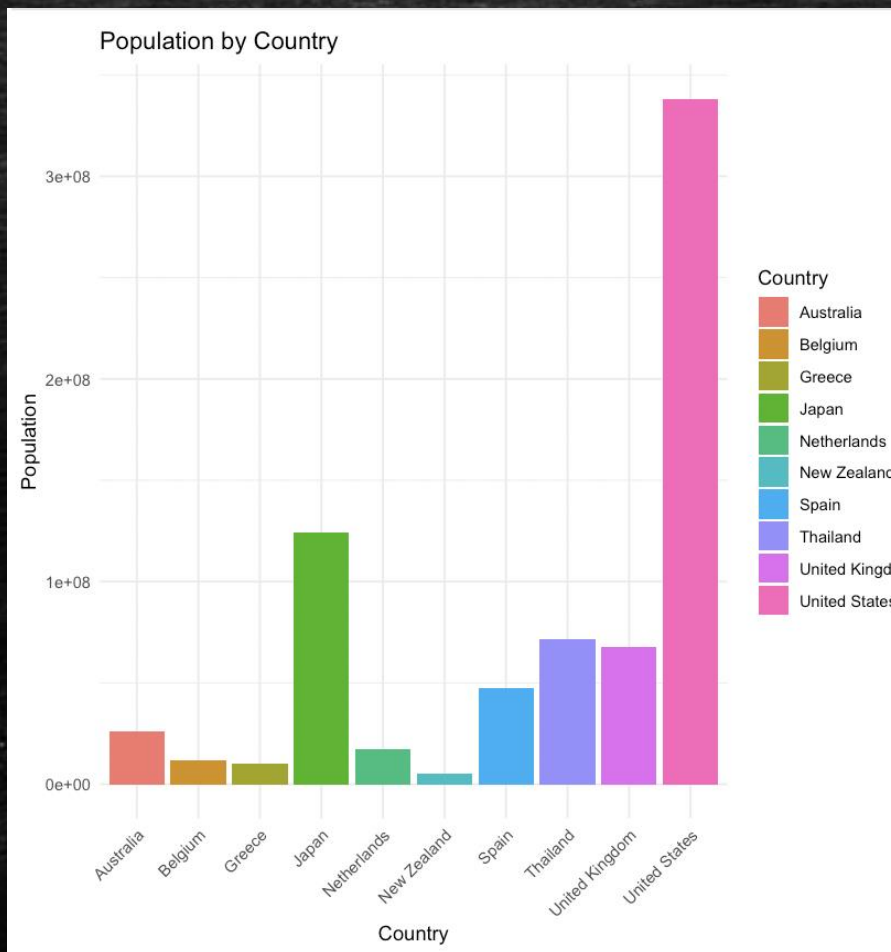
japan_airbnb_nzd
```

A tibble: 12234 × 3

price_nzd	bathrooms_text	bedrooms
<dbl>	<chr>	<dbl>
133.77839	1 bath	1
103.27691	1 bath	NA
82.70849	1 bath	1
99.24127	1 shared bath	NA

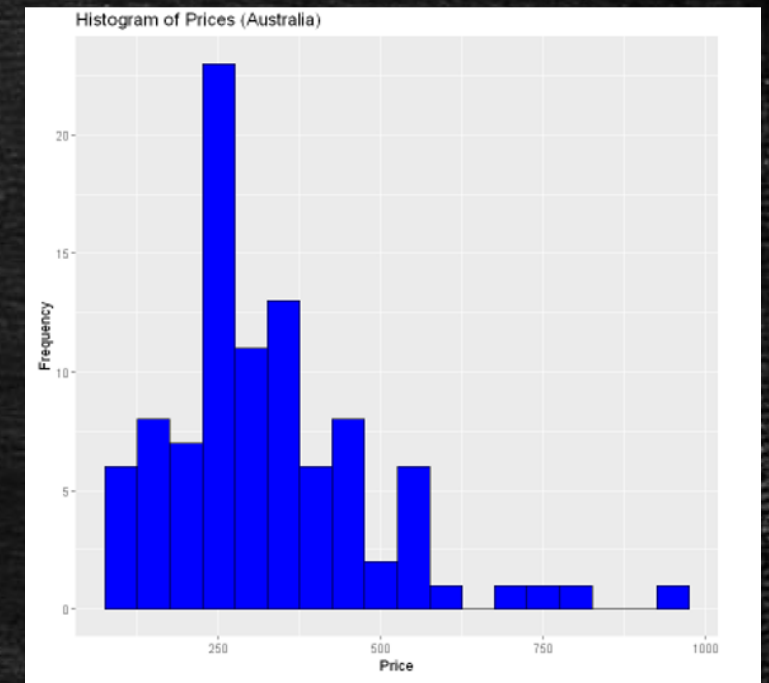
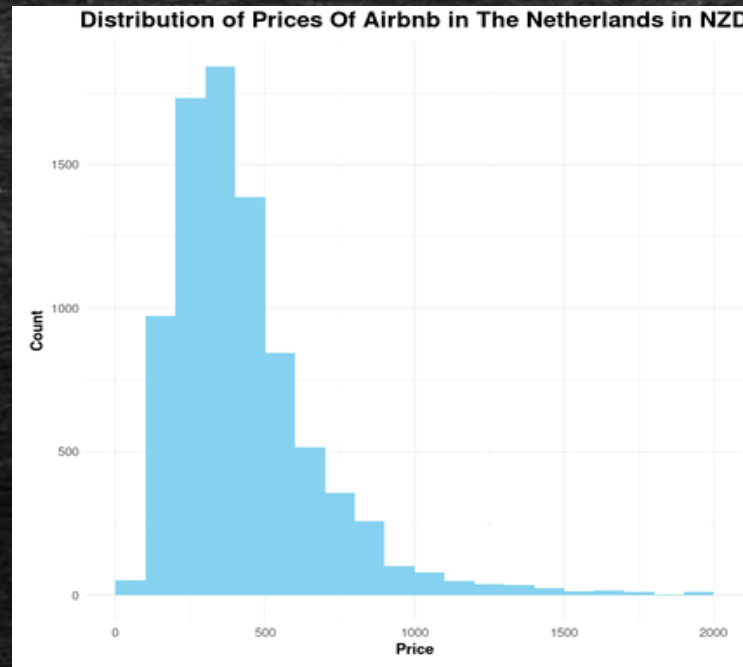
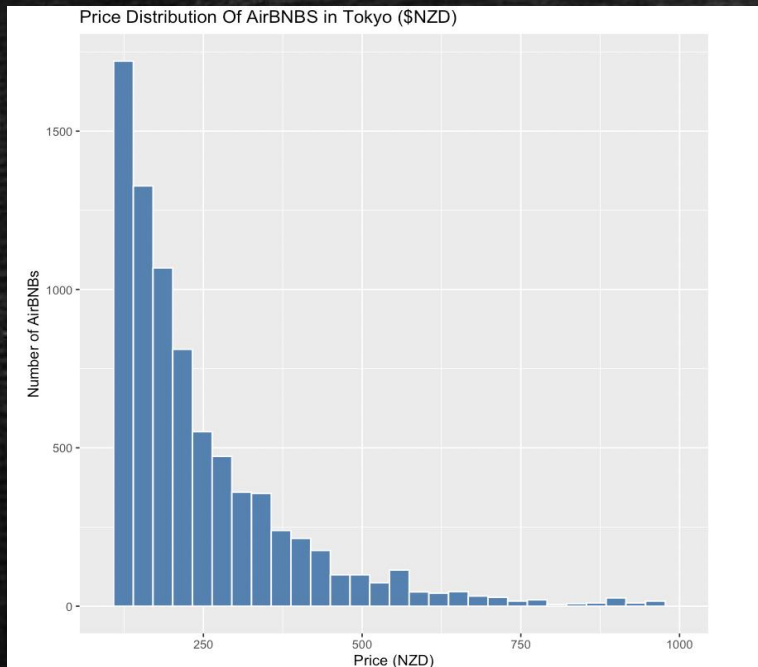
# Our findings!

Before we look into the actual AirBNB data here are the graphs for the population and GDP per capita for each of the countries we scraped AirBNB data for



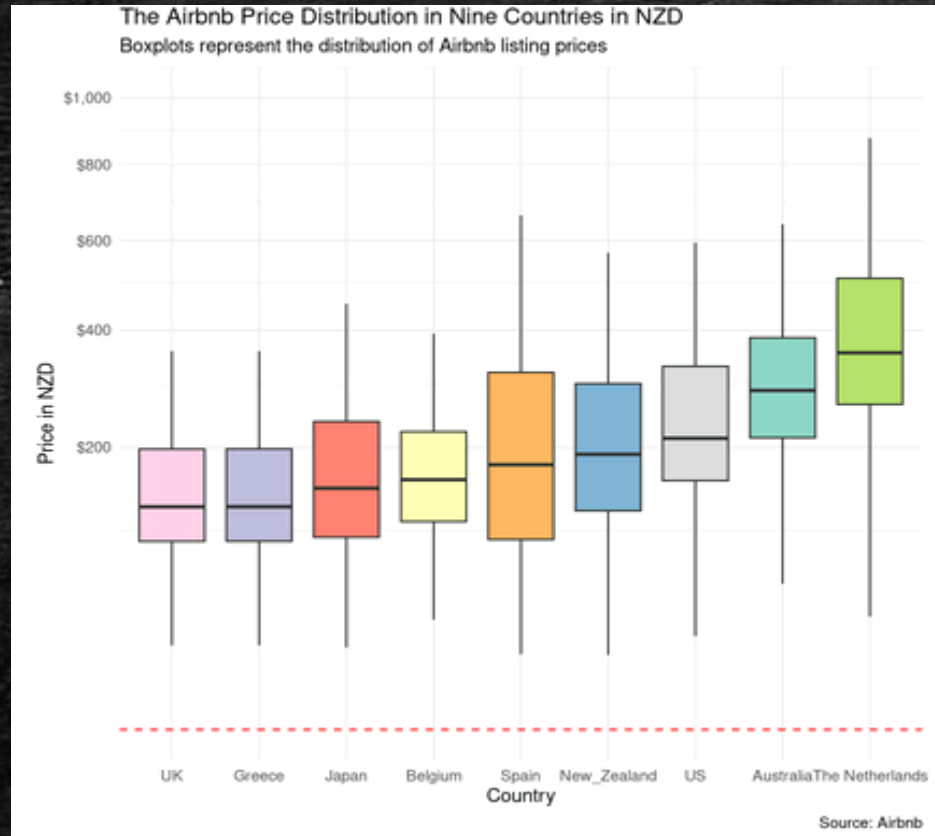


# AirBNB Price Distribution Graphs



In these graphs we can see that the majority of AirBNB prices are at the lower end being between \$0-\$250, and the number of listings significantly drop at price points above \$500. When comparing these graphs to previous graphs of population and GDP per capita there is no direct comparison, that could help us answer our research question.

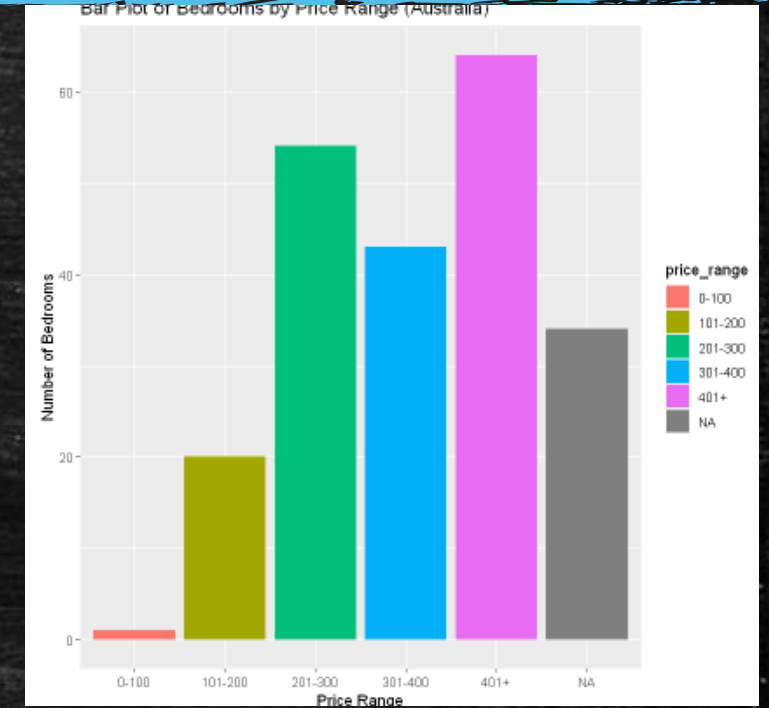
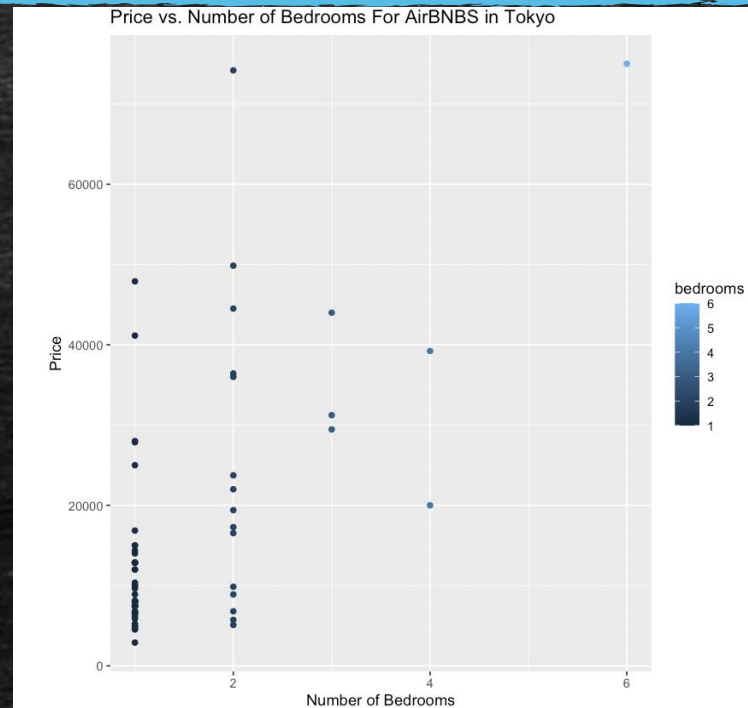
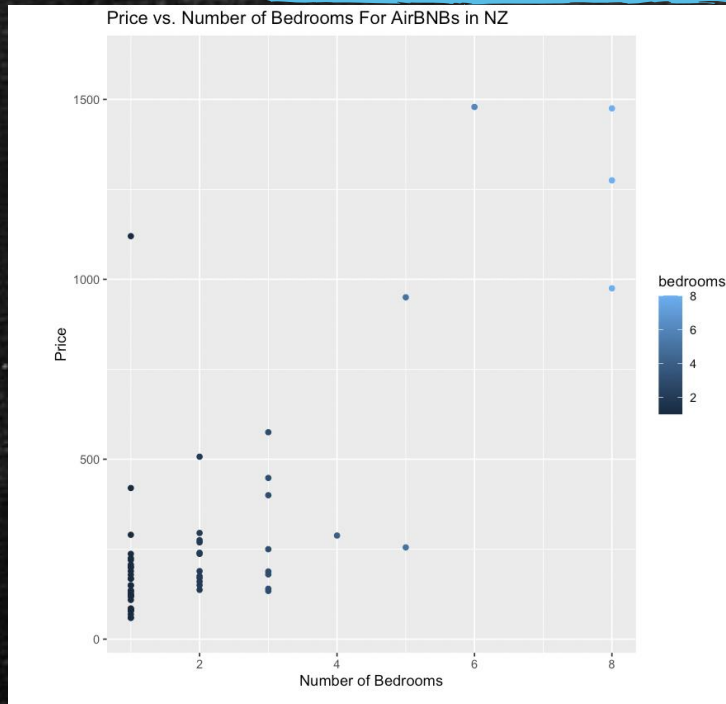
# BOX PLOT OF PRICE DISTRIBUTION IN THREE COUNTRIES



This graph shows us that The Netherlands, Australia and USA have the highest price distribution of AirBNBs. These three countries also have the highest GDP per capita. This graph also shows us that Greece has one of lowest price distributions for AirBNBs, and Greece had the lowest GDP per capita between these countries. This answers a part of our research question where the higher the GDP per capita the higher the price of the AirBNBs and vice versa. However, there was no relation to the country's population and AirBNB pricing

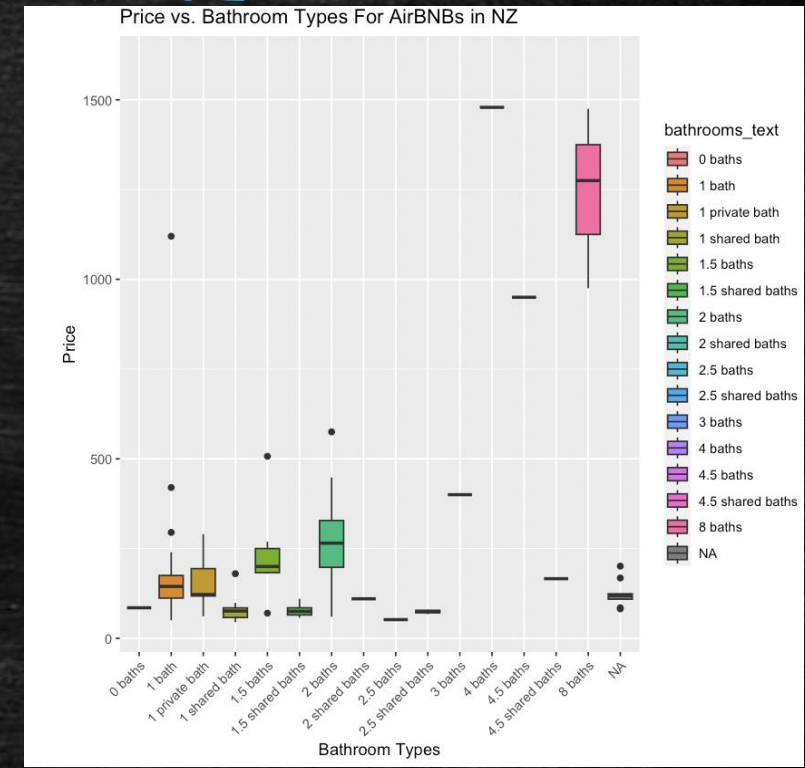
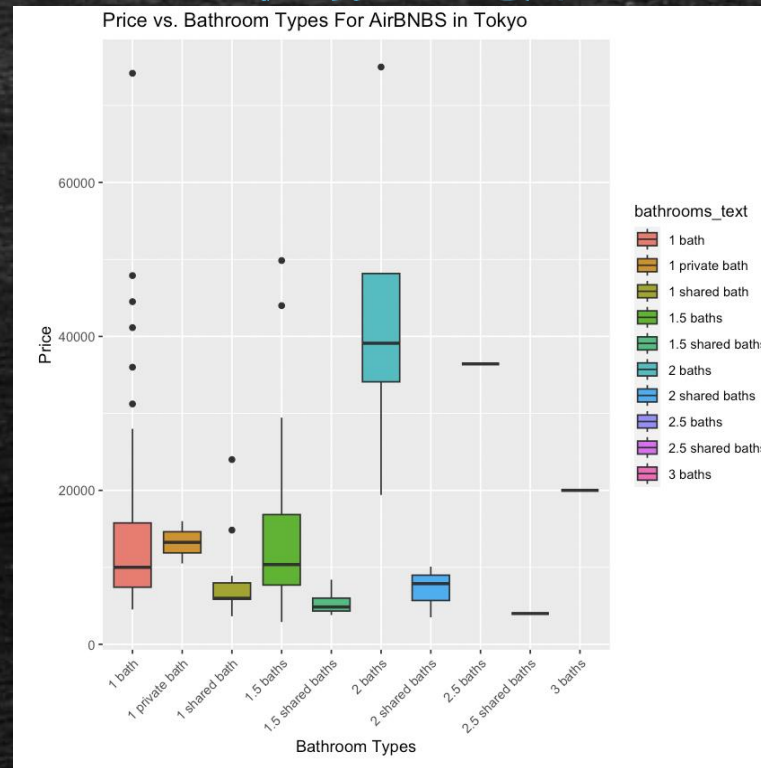
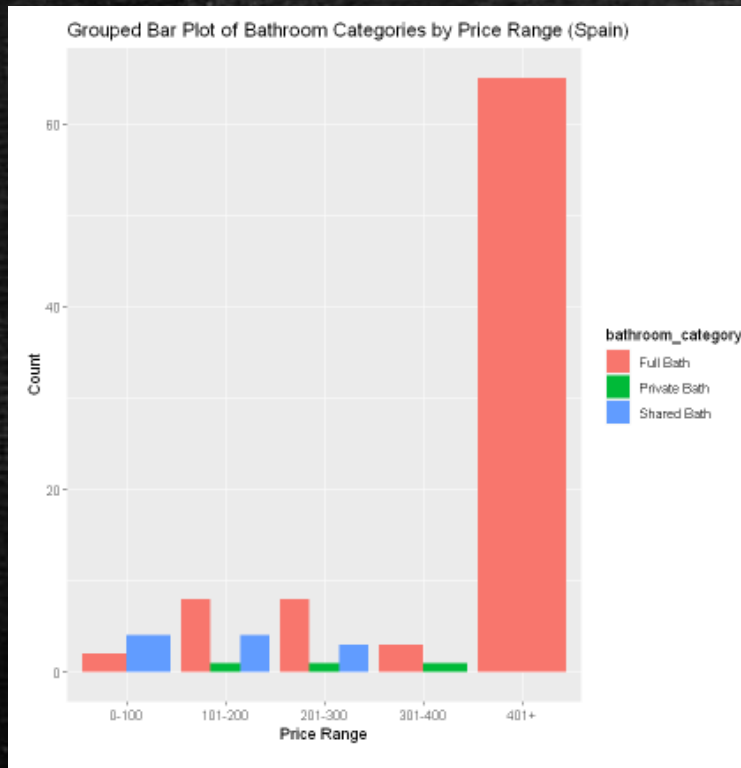


# Comparing AirBNB Prices With The Number Of Bedrooms



These graphs show us that the more bedrooms AirBNBs have the higher priced they tend to be. For example, the highest priced AirBNB in Tokyo has the most bedrooms which is 6.

# Comparing AirBNB Prices With The Number and Type Of Bathrooms



These graphs shows us that the more bathrooms the higher the AirBNB Price tends to be especially the price increase between 1-2 bathrooms. However the price increase when there was more than 2 bathrooms wasn't as great as the price difference between 1 and 2 bathrooms



# Conclusion

---

- Our research question was if population and GDP per capita had any correlation with the pricing of AirBNBs. However, we could not conclude that population has any effects on AirBNB pricing.
- GDP per capita did have correlation with the prices of AirBNBs as seen in slide 10 and 12, where the 3 countries with the highest GDP per capita had the highest price distributions of AirBNBs and the countries with the lowest GDP per capita had the lowest price distribution of AirBNBs.
- The second part of our question was if the AirBNBs had a higher bedrooms and bathrooms does that mean the pricing would be higher, and yes from our findings we can conclude that this was true.

Thanks for listening!

---



---

**Q & A**