

# Assignment 2 Group 8

Due at 11:59pm on October 3.

Leng Seong Che; Bozhou (Peter) Tan

The results are based on the data at 2:40pm on Oct 4th.

GitHub repo: <https://github.com/petertbz/Assignment-2--727.git>

!! Notice: The data from Google API is changing every second, so our analysis is based on the data we obtained at the time we mention above.

You may work in pairs or individually for this assignment. Make sure you join a group in Canvas if you are working in pairs. Turn in this assignment as an HTML or PDF file to ELMS. Make sure to include the R Markdown or Quarto file that was used to generate it.

```
library(tidyverse)
library(gtrendsR)
library(censusapi)
```

In this assignment, you will pull from APIs to get data from various data sources and use your data wrangling skills to use them all together. You should turn in a report in PDF or HTML format that addresses all of the questions in this assignment, and describes the data that you pulled and analyzed. You do not need to include full introduction and conclusion sections like a full report, but you should make sure to answer the questions in paragraph form, and include all relevant tables and graphics.

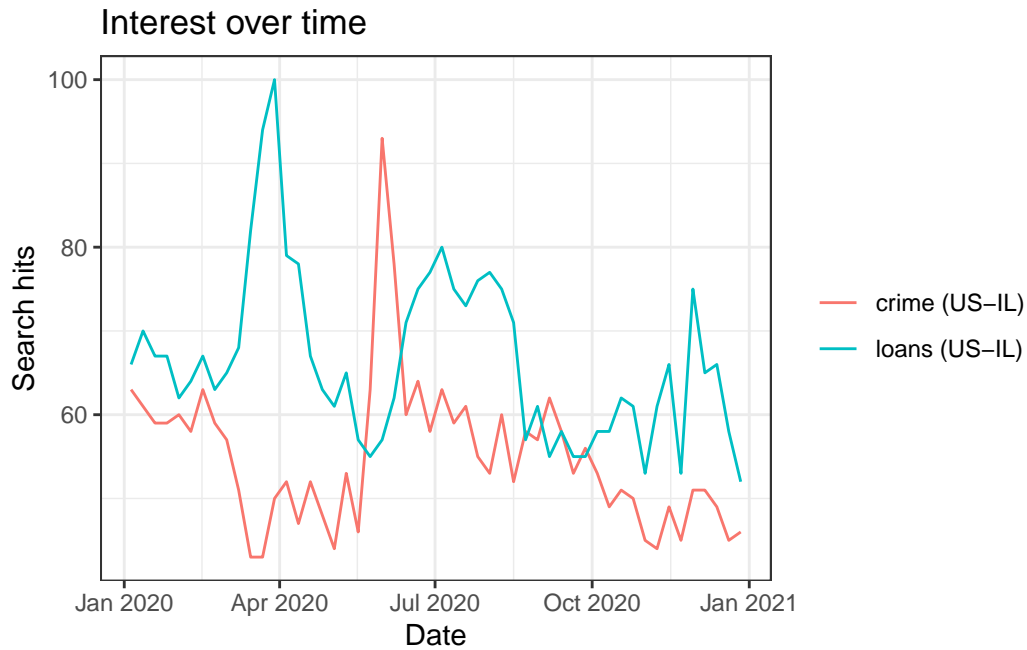
Whenever possible, use piping and `dplyr`. Avoid hard-coding any numbers within the report as much as possible.

## Pulling from APIs

### crime and loans

Our first data source is the Google Trends API. Suppose we are interested in the search trends for `crime` and `loans` in Illinois in the year 2020. We could find this using the following code:

```
res = gtrends(c("crime", "loans"),
              geo = "US-IL",
              time = "2020-01-01 2020-12-31",
              low_search_volume = TRUE)
plot(res)
```



Answer the following questions for the keywords “crime” and “loans”.

- Find the mean, median and variance of the search hits for the keywords.

```
# transfer the data into tibble
rest = as_tibble(res$interest_over_time)

# find the mean, median and variance of the search hits
library(dplyr)
library(tidyr)
library(knitr)

descriptive = rest %>%
  group_by(keyword) %>%
  summarise(n = n(),
```

```

mean = mean(hits),
median = median(hits),
variance = var(hits))
kable(descriptive, caption = "Descriptive Statistics of Keywords")

```

Table 1: Descriptive Statistics of Keywords

keyword	n	mean	median	variance
crime	52	54.98077	53	78.13688
loans	52	66.50000	65	101.39216

According to Table 1, we can find that the keyword crime has a mean of 54.9807692307692, a median of 53 and a variance of 78.1368778280543. The keyword loans has a mean of 66.5, a median of 65 and a variance of 101.392156862745.

- Which cities (locations) have the highest search frequency for `loans`? Note that there might be multiple rows for each city if there were hits for both “crime” and “loans” in that city. It might be easier to answer this question if we had the search hits info for both search terms in two separate variables. That is, each row would represent a unique city.

```

rescity = as_tibble(res$interest_by_city) %>%
  pivot_wider(., names_from = keyword, values_from = hits) %>%
  arrange(., desc(loans))
kable(head(rescity), caption = "Highest Search Frequency for Loans")

```

Table 2: Highest Search Frequency for Loans

location	geo	gprop	crime	loans
Hinckley	US-IL	web	NA	100
Carrier Mills	US-IL	web	NA	96
Glasford	US-IL	web	NA	94
Riverton	US-IL	web	NA	88
Georgetown	US-IL	web	NA	88
Rosemont	US-IL	web	44	87

According to Table 2, Hinckley has the highest search frequency for `loans` with the value of 100, followed by Carrier Mills and Glasford.

- Is there a relationship between the search intensities between the two keywords we used?

```

crime = rest %>%
  filter(keyword == "crime") %>%
  select(date, hits) %>%
  rename(., crimehits = hits)

loan = rest %>%
  filter(keyword == "loans") %>%
  select(date, hits) %>%
  rename(., loanshits = hits)

crimloan = left_join(crime, loan, by = "date")
cor.test(crimloan$crimehits, crimloan$loanshits)

```

Pearson's product-moment correlation

```

data: crimloan$crimehits and crimloan$loanshits
t = -0.62945, df = 50, p-value = 0.5319
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 -0.3530257  0.1888011
sample estimates:
      cor
-0.08866771

```

According to the plot at the beginning, search frequencies for “crime” and “loans” have similar trends at the beginning of 2020, where they both went up and down from January to around February 2020. From March to April, search frequency for “loans” increased drastically from approximately 65 to 100, while search frequency for “crime” decreased before it increased again. In other words, the two keywords have a similar trend between January and February and most time between July 2020 and January 2021. However, from March to June 2020, they seem to have an inverse relationship.

If we use the quantitative method to compute the t-statistic and corresponding p-value, we can see that the p-value is bigger than 0.05, which means there is no statistically significant negative relationship between crime and loans at significance level 0.05.

```
cor.test(rescity$crime, rescity$loans)
```

Pearson's product-moment correlation

```
data: rescity$crime and rescity$loans
t = 0.49472, df = 14, p-value = 0.6285
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 -0.3899644  0.5885433
sample estimates:
      cor
0.1310796
```

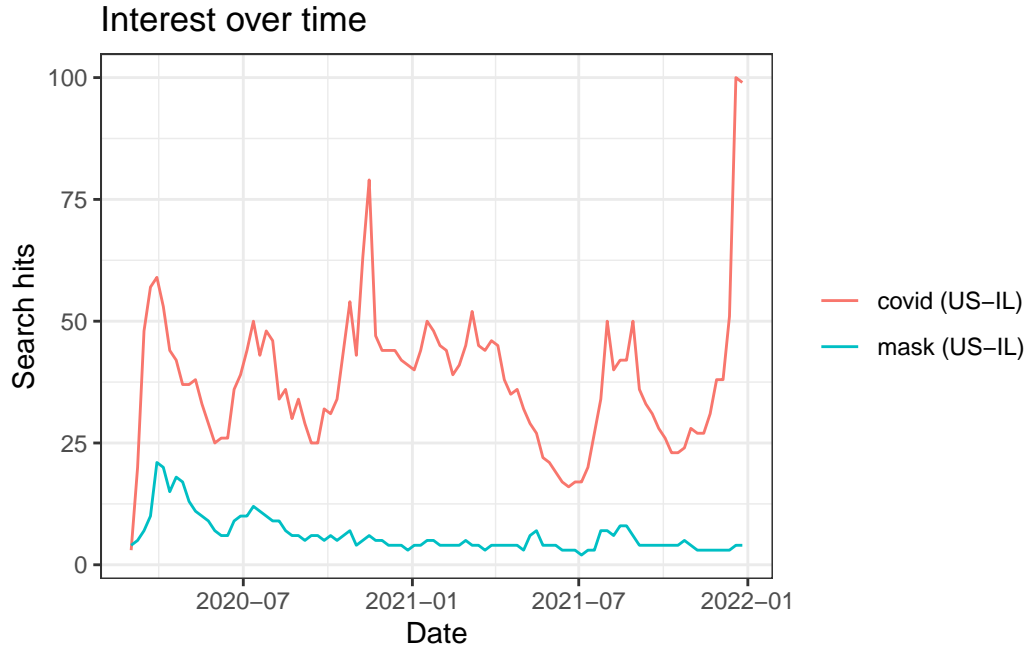
If we investigate the relationship based on the interest by city data, the correlation test result suggests that search frequencies for “crime” and “loans” do not have a statistically significant correlation. However, lots of data are missing in this dataset, which might affect the reliability of the correlation test result.

### **covid and mask**

Repeat the above for keywords related to covid. Make sure you use multiple keywords like we did above. Try several different combinations and think carefully about words that might make sense within this context.

We choose **covid** and **mask** as our keywords for analysis.

```
res2 = gtrends(c("covid", "mask"),
               geo = "US-IL",
               time = "2020-03-01 2021-12-31",
               low_search_volume = TRUE)
plot(res2)
```



```
# transfer the data into tibble
rest2 = as_tibble(res2$interest_over_time)

# find the mean, median and variance of the search hits
descriptive2 = rest2 %>%
  group_by(keyword) %>%
  summarise(n = n(),
            mean = mean(hits),
            median = median(hits),
            variance = var(hits))
kable(descriptive2, caption = "Descriptive Statistics of Keywords")
```

Table 3: Descriptive Statistics of Keywords

keyword	n	mean	median	variance
covid	96	38.125000	38	218.80526
mask	96	6.083333	5	13.78246

From the table, we can find that the keyword covid has a mean of 38.125, a median of 38 and a variance of 218.805263157895. The keyword mask has a mean of 6.08333333333333, a median of 5 and a variance of 13.7824561403509.

```

rescity2 = as_tibble(res2$interest_by_city) %>%
  pivot_wider(., names_from = keyword, values_from = hits) %>%
  arrange(., desc(covid))
kable(head(rescity2), caption = "Highest Search Frequency for covid")

```

Table 4: Highest Search Frequency for covid

location	geo	gprop	covid	mask
Hinsdale	US-IL	web	100	NA
Barrington	US-IL	web	98	NA
Naperville	US-IL	web	95	NA
Highland Park	US-IL	web	94	80
Northbrook	US-IL	web	94	NA
Oak Lawn	US-IL	web	93	NA

From the table, we can see that Hinsdale has the highest search frequency for covid with the value of 100, followed by Barrington and Naperville.

```

mask = rest2 %>%
  filter(keyword == "mask") %>%
  select(date, hits) %>%
  rename(., maskhits = hits)

covid = rest2 %>%
  filter(keyword == "covid") %>%
  select(date, hits) %>%
  rename(., covidhits = hits)

maskcovid = left_join(mask, covid, by = "date")
cor.test(maskcovid$maskhits, maskcovid$covidhits)

```

Pearson's product-moment correlation

```

data: maskcovid$maskhits and maskcovid$covidhits
t = 2.1751, df = 94, p-value = 0.03213
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.01926168 0.40175642
sample estimates:

```

```
cor
0.2189023
```

From the correlation test, we can see that covid has a significantly positive correlation with mask at 0.05 level. The correlation probably means that people will search for mask when Covid-19 is severe in one place.

```
cor.test(rescity2$covid, rescity2$mask)
```

Pearson's product-moment correlation

```
data: rescity2$covid and rescity2$mask
t = 3.9658, df = 44, p-value = 0.0002656
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.2618548 0.6992759
sample estimates:
      cor
0.5131456
```

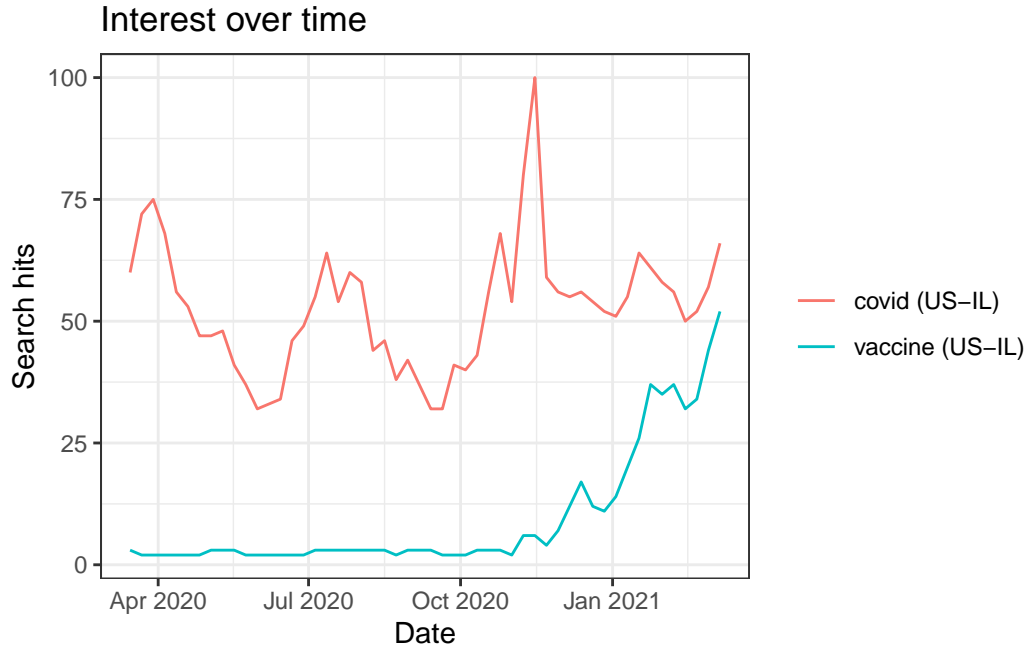
If we investigate the relationship based on the interest by city data, the correlation test result suggests a significantly positive relationship between search frequencies for covid and mask at 0.05 significance level.

## covid and vaccine

We choose covid and vaccine as our keywords for analysis.

```
res3 = gtrends(c("covid", "vaccine"),
               geo = "US-IL",
               time = "2020-03-11 2021-3-11",
               low_search_volume = TRUE)
plot(res3)
```





```
# transfer the data into tibble
rest3 = as_tibble(res3$interest_over_time)

# find the mean, median and variance of the search hits
descriptive3 = rest3 %>%
  group_by(keyword) %>%
  summarise(n = n(),
            mean = mean(hits),
            median = median(hits),
            variance = var(hits))
kable(descriptive3, caption = "Descriptive Statistics of Keywords")
```

Table 5: Descriptive Statistics of Keywords

keyword	n	mean	median	variance
covid	52	52.769231	54	171.9849
vaccine	52	9.442308	3	166.1731

From the table, we can find that the keyword covid has a mean of 52.7692307692308, a median of 54 and a variance of 171.984917043741. The keyword vaccine has a mean of 9.44230769230769, a median of 3 and a variance of 166.173076923077.

```

rescity3 = as_tibble(res3$interest_by_city) %>%
  pivot_wider(., names_from = keyword, values_from = hits) %>%
  arrange(., desc(vaccine))
kable(head(rescity3), caption = "Highest Search Frequency for vaccine")

```

Table 6: Highest Search Frequency for vaccine

location	geo	gprop	covid	vaccine
Mount Sterling	US-IL	web	NULL	15
Villa Grove	US-IL	web	NULL	16
Watseka	US-IL	web	NULL	23
Toulon	US-IL	web	NULL	23
Philo	US-IL	web	NULL	26
Jerseyville	US-IL	web	NULL	26

From the table, we can see that Mount Sterling has the highest search frequency for vaccine with the value of 100, followed by Villa Grove and Watseka.

```

vaccine = rest3 %>%
  filter(keyword == "vaccine") %>%
  select(date, hits) %>%
  rename(., vaccinehits = hits)

covid = rest3 %>%
  filter(keyword == "covid") %>%
  select(date, hits) %>%
  rename(., covidhits = hits)

vaccinecovid = left_join(vaccine, covid, by = "date")
cor.test(vaccinecovid$vaccinehits, vaccinecovid$covidhits)

```

Pearson's product-moment correlation

```

data: vaccinecovid$vaccinehits and vaccinecovid$covidhits
t = 1.6456, df = 50, p-value = 0.1061
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
-0.04927685 0.47046934
sample estimates:

```

```
cor
0.2266717
```

The search popularities for covid and vaccine seem to have a positive relationship starting from November 2020. As the vaccine became available around this time, the search for vaccine increases drastically and follows a similar pattern of covid.

However, the correlation test suggests that no statistically significant relationship appears between covid and vaccine at 0.05 significance level.

```
rescity3$covid <- as.numeric(as.character(rescity3[[4]]))
rescity3$vaccine <- as.numeric(as.character(rescity3[[5]]))
cor.test(rescity3$covid, rescity3$vaccine)
```

Pearson's product-moment correlation

```
data: rescity3$covid and rescity3$vaccine
t = 6.5202, df = 48, p-value = 4.005e-08
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.5029306 0.8093016
sample estimates:
      cor
0.6853382
```

If we investigate the relationship based on the interest by city data, the correlation test result is inconsistent with the one based on the over time data. This one suggests a significantly positive relationship between search frequencies for covid and vaccine at 0.05 significance level.

## Google Trends + ACS

### crime and loans

Now lets add another data set. The `censusapi` package provides a nice R interface for communicating with this API. However, before running queries we need an access key. This (easy) process can be completed here:

[https://api.census.gov/data/key\\_signup.html](https://api.census.gov/data/key_signup.html)

Once you have an access key, store this key in the `cs_key` object. We will use this object in all following API queries.

```
cs_key <- "c0fd12402e23b7a95923e694f046015d624c91c5"
```

In the following, we request basic socio-demographic information (population, median age, median household income, income per capita) for cities and villages in the state of Illinois.

```
acs_il <- getCensus(name = "acs/acs5",
  vintage = 2020,
  vars = c("NAME",
    "B01001_001E",
    "B06002_001E",
    "B19013_001E",
    "B19301_001E"),
  region = "place:*",
  regionin = "state:17",
  key = cs_key)

head(acs_il)
```

	state	place	NAME	B01001_001E	B06002_001E	B19013_001E
1	17	15261 Coatsburg village, Illinois		180	35.6	55714
2	17	15300 Cobden village, Illinois		1018	44.2	38750
3	17	15352 Coffeen city, Illinois		640	33.4	35781
4	17	15378 Colchester city, Illinois		1347	42.2	43942
5	17	15469 Coleta village, Illinois		230	27.7	56875
6	17	15495 Colfax village, Illinois		1088	32.5	58889
		B19301_001E				
1		27821				
2		19979				
3		26697				
4		24095				
5		23749				
6		24861				

Convert values that represent missings to NAs.

```
acs_il[acs_il == -666666666] <- NA
```

Now, it might be useful to rename the socio-demographic variables (B01001\_001E etc.) in our data set and assign more meaningful names.

```
acs_il <-
  acs_il %>%
```

```
rename(pop = B01001_001E,
       age = B06002_001E,
       hh_income = B19013_001E,
       income = B19301_001E)
```

It seems like we could try to use this location information listed above to merge this data set with the Google Trends data. However, we first have to clean `NAME` so that it has the same structure as `location` in the search interest by city data. Add a new variable `location` to the ACS data that only includes city names.

```
library(stringr)
pattern = c("St." = "Saint")

acs_il = acs_il %>%
  mutate(location = str_remove_all(NAME, c(" town,| city,| village,| Illinois"))) %>%
  mutate(location = str_replace_all(location, coll(pattern)))
```

Answer the following questions with the “crime” and “loans” Google trends data and the ACS data.

- First, check how many cities don’t appear in both data sets, i.e. cannot be matched. Then, create a new data set by joining the Google Trends and the ACS data. Keep only cities that appear in both data sets.

```
joint = inner_join(rescity, acs_il, by = "location")
nrow(joint)
```

[1] 335

```
# check how many cities do not appear in both datasets
n = (nrow(acs_il) - nrow(joint) ) + (nrow(rescity) - nrow(joint))
n
```

[1] 1142

1142 cities do not appear in both datasets.

- Compute the mean of the search popularity for both keywords for cities that have an above average median household income and for those that have an below average median household income. When building your pipe, start with creating the grouping variable and then proceed with the remaining tasks. What conclusions might you draw from this?

```
group1 = joint %>%
  mutate(mean = mean(hh_income, na.rm = TRUE)) %>%
  mutate(group = ifelse(hh_income > mean, "high", "low")) %>%
  group_by(group) %>%
  summarise(crime = mean(crime, na.rm = TRUE),
            loans = mean(loans, na.rm = TRUE)) %>%
  filter(!is.na(group))
kable(group1, caption = "Search Popularity by Household Income")
```

Table 7: Search Popularity by Household Income

group	crime	loans
high	45.26316	46.94643
low	50.81250	52.97531

From the table, cities that have an above average median household income have lower crime hits and lower loans hits, which means crime and loans may correlate with income. The reason for higher mean search popularity of “crime” can be that those with below average median household income live in some neighborhoods with a relatively higher number of crimes. Houses in areas with more crimes can be more affordable. The reason for higher mean search popularity of “loans” can be these households need more loans for various living expenses such as education.

- Is there a relationship between the median household income and the search popularity of the Google trends terms? Describe the relationship and use a scatterplot with `qplot()`.

```
cor.test(joint$hh_income, joint$crime, method = "pearson")
```

Pearson's product-moment correlation

```
data: joint$hh_income and joint$crime
t = -1.8423, df = 68, p-value = 0.06979
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 -0.43093153  0.01785344
sample estimates:
cor
-0.2180353
```

```
cor.test(joint$hh_income, joint$loans, method = "pearson")
```

Pearson's product-moment correlation

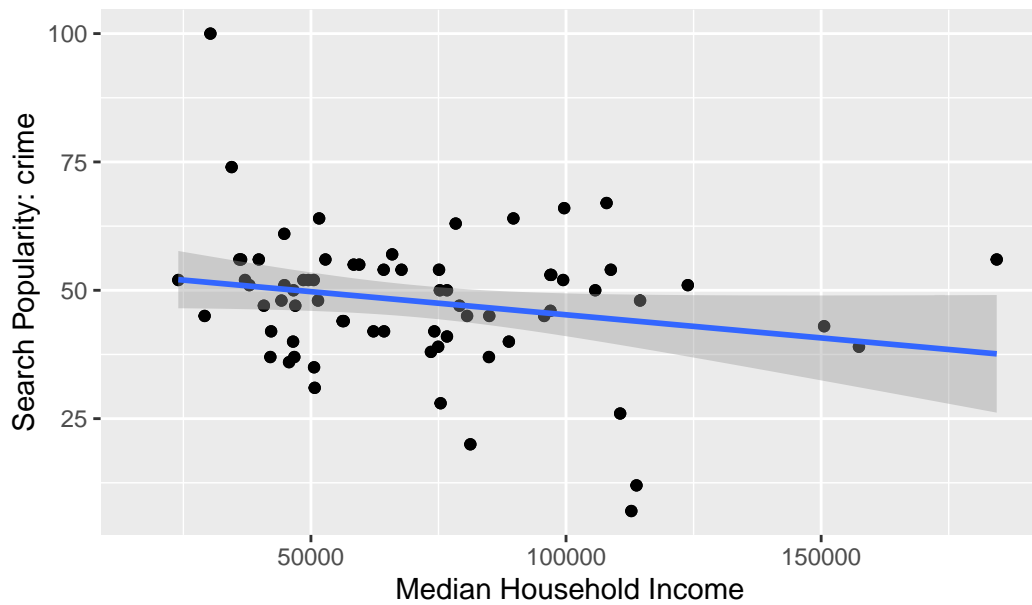
```
data: joint$hh_income and joint$loans
t = -2.1798, df = 135, p-value = 0.03101
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 -0.3415443 -0.0172078
sample estimates:
      cor
-0.1843911
```

```
p1 = qplot(x = hh_income, y = crime, data = joint) +
  geom_point() +
  geom_smooth(method = lm) +
  labs(
    title = "Scatter Plot of Median Household Income vs. 'crime' Search by City",
    x = "Median Household Income",
    y = "Search Popularity: crime"
  )

p1
```

`geom\_smooth()` using formula = 'y ~ x'

Scatter Plot of Median Household Income vs. 'crime' Search by

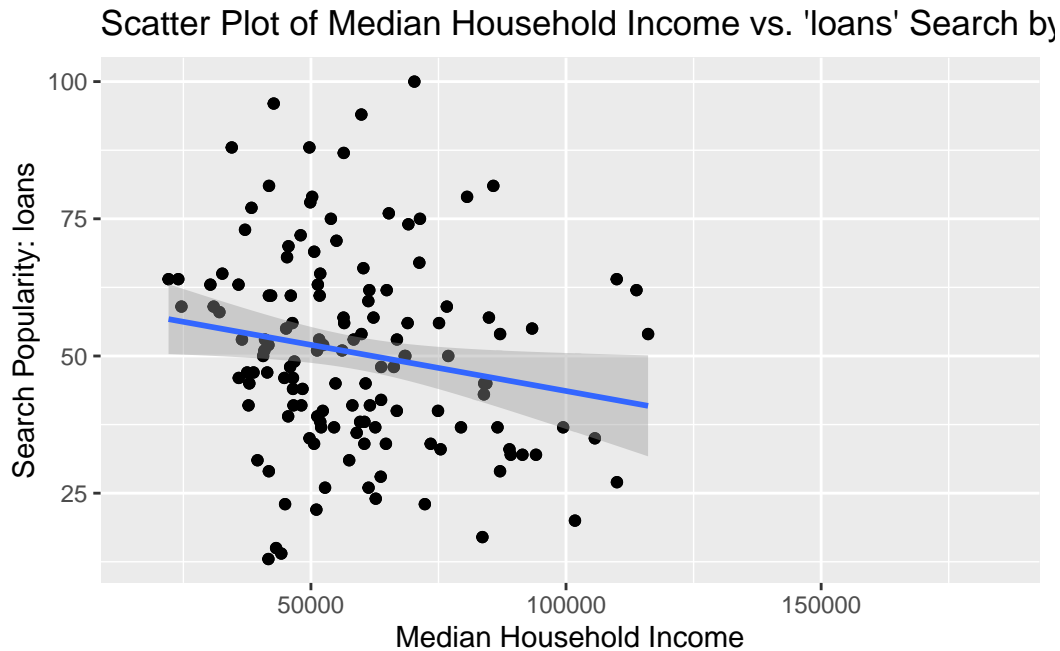


```
p2 = qplot(x = hh_income, y = loans, data = joint) +
  geom_point() +
  geom_smooth(method = lm) +
  labs(
    title = "Scatter Plot of Median Household Income vs. 'loans' Search by City",
    x = "Median Household Income",
    y = "Search Popularity: loans"
  )
```

p2

`geom\_smooth()` using formula = 'y ~ x'





The results from the Pearson correlation test suggest a negative statistically significant correlation between the median household income and the search popularity for “loans” and a statistically non-significant correlation between the median household income and “crime” at 0.05 level. These can also be observed in the scatter plots. For “crime”, the majority of the cities have search popularity below 60 regardless of median household income. A slightly decreasing trend according to the the regression line, however, the correlation test suggests an absence of statistically significant relationship between them. For “loans”, a decreasing trend is suggested based on the regression line. As the median household income increases, the search popularity for “loans” decrease. For those with median household income higher than \$80,000, the searches are mostly lower than 50. In contrast, those with median household income lower than \$80,000 have a wider range of search numbers, and half cities among this group have over 50 loans hits.

### covid and mask

Repeat the above steps using the covid data and the ACS data.

```
joint2 = inner_join(rescity2, acs_il, by = "location")

# check how many cities do not appear in both datasets
n2 = (nrow(acs_il) - nrow(joint2) ) + (nrow(rescity2) - nrow(joint2))
```

n2

[1] 1137

```
group2 = joint2 %>%
  mutate(mean = mean(hh_income, na.rm = TRUE)) %>%
  mutate(group = ifelse(hh_income > mean, "high", "low")) %>%
  group_by(group) %>%
  summarise(covid = mean(covid, na.rm = TRUE),
            mask = mean(mask, na.rm = TRUE)) %>%
  filter(!is.na(group))
kable(group2, caption = "Search Popularity by Household Income")
```

Table 8: Search Popularity by Household Income

group	covid	mask
high	73.78667	70.04615
low	59.52137	66.02941

1137 cities do not appear in both datasets.

From the table, we can see cities that have an above average median household income have higher covid hits and higher mask hits, which means search hits of covid and mask may correlate with income positively.

```
cor.test(joint2$hh_income, joint2$covid, method = "pearson")
```

Pearson's product-moment correlation

```
data: joint2$hh_income and joint2$covid
t = 12.301, df = 190, p-value < 2.2e-16
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.5788060 0.7378749
sample estimates:
      cor
0.6658396
```

```
cor.test(joint2$hh_income, joint2$mask, method = "pearson")
```

Pearson's product-moment correlation

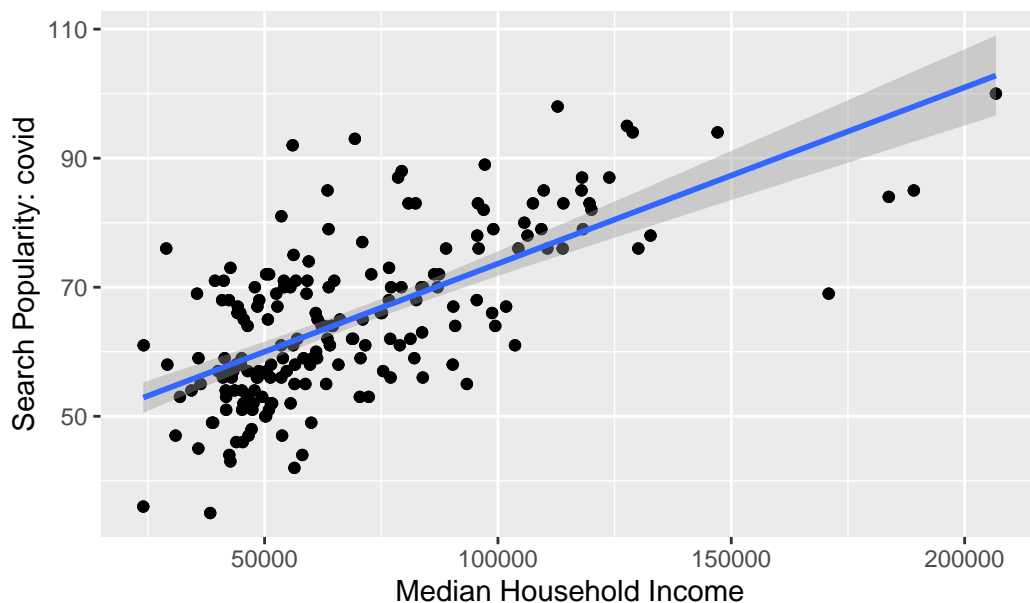
```
data: joint2$hh_income and joint2$mask
t = 2.6134, df = 165, p-value = 0.009793
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.04898821 0.34091051
sample estimates:
      cor
0.1993686
```

```
p3 = qplot(x = hh_income, y = covid, data = joint2) +
  geom_point() +
  geom_smooth(method = lm) +
  labs(
    title = "Scatter Plot of Median Household Income vs. 'covid' Search by City",
    x = "Median Household Income",
    y = "Search Popularity: covid"
  )

p3
```

`geom\_smooth()` using formula = 'y ~ x'

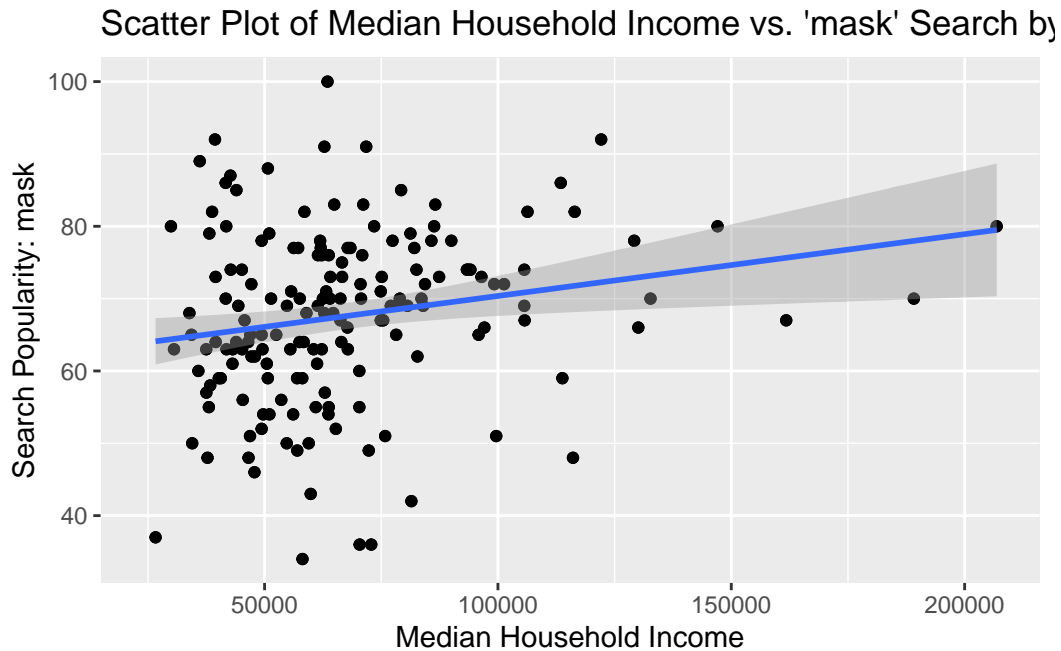
Scatter Plot of Median Household Income vs. 'covid' Search by



```
p4 = qplot(x = hh_income, y = mask, data = joint2) +
  geom_point() +
  geom_smooth(method = lm) +
  labs(
    title = "Scatter Plot of Median Household Income vs. 'mask' Search by City",
    x = "Median Household Income",
    y = "Search Popularity: mask"
  )
```

p4

`geom\_smooth()` using formula = 'y ~ x'



According to the scatterplots, we can see that the income is both positively correlated with covid and mask. This indicates that people in rich areas may pay attention to covid and its protection more, showing one kind of social inequality. From the correlation test, we can see that the p-value of both tests are both less than 0.05, indicating that income has a statistically significant positive relation with both covid and mask.

### covid and vaccine

```
joint3 = inner_join(rescity3, acs_il, by = "location")

# check how many cities do not appear in both datasets
n3 = (nrow(acs_il) - nrow(joint3)) + (nrow(rescity3) - nrow(joint3))
n3
```

[1] 1148

```
group3 = joint3 %>%
  mutate(mean = mean(hh_income, na.rm = TRUE)) %>%
  mutate(group = ifelse(hh_income > mean, "high", "low")) %>%
  group_by(group) %>%
```

```

summarise(covid = mean(covid, na.rm = TRUE),
          vaccine = mean(vaccine, na.rm = TRUE)) %>%
filter(!is.na(group))
kable(group3, caption = "Search Popularity by Household Income")

```

Table 9: Search Popularity by Household Income

group	covid	vaccine
high	77.46429	53.14943
low	65.74400	35.62069

1148 cities do not appear in both datasets.

From the table, we can see cities that have an above average median household income have higher covid hits and higher mask hits, which means search hits of covid and mask may correlate with income positively.

```

cor.test(joint3$hh_income, joint3$covid, method = "pearson")

```

Pearson's product-moment correlation

```

data: joint3$hh_income and joint3$covid
t = 8.7196, df = 179, p-value = 1.873e-15
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.4347766 0.6408311
sample estimates:
      cor
0.5460089

```

```

cor.test(joint3$hh_income, joint3$vaccine, method = "pearson")

```

Pearson's product-moment correlation

```

data: joint3$hh_income and joint3$vaccine
t = 11.514, df = 172, p-value < 2.2e-16
alternative hypothesis: true correlation is not equal to 0

```

95 percent confidence interval:

0.5666078 0.7362666

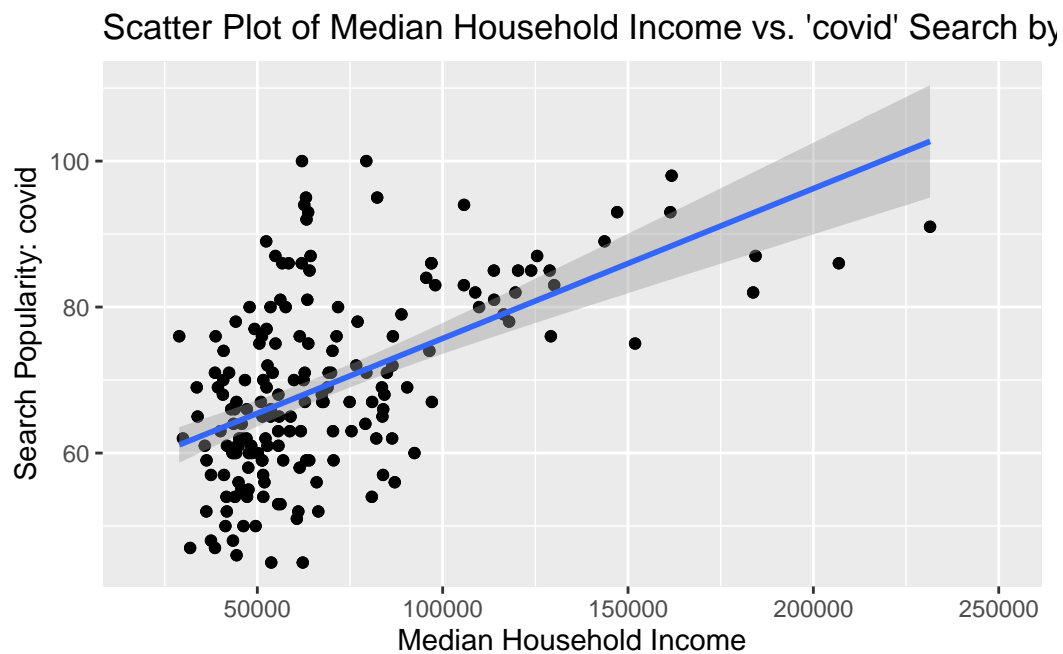
sample estimates:

cor

0.6597634

```
p5 = qplot(x = hh_income, y = covid, data = joint3) +  
  geom_point() +  
  geom_smooth(method = lm) +  
  labs(  
    title = "Scatter Plot of Median Household Income vs. 'covid' Search by City",  
    x = "Median Household Income",  
    y = "Search Popularity: covid"  
  )  
  
p5
```

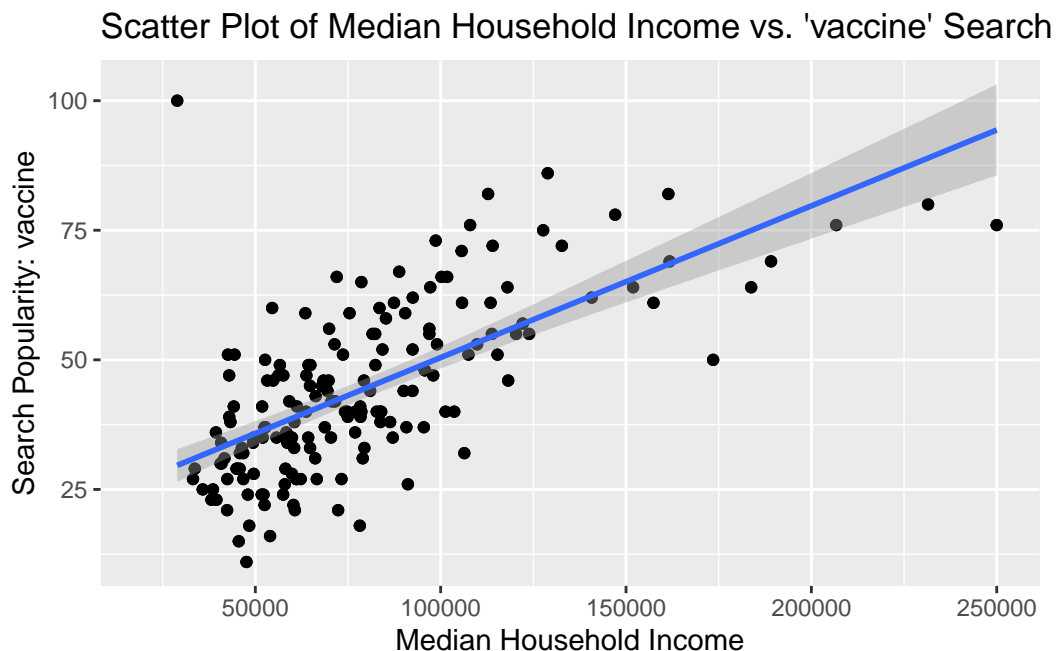
`geom\_smooth()` using formula = 'y ~ x'



```
p6 = qplot(x = hh_income, y = vaccine, data = joint3) +
  geom_point() +
  geom_smooth(method = lm) +
  labs(
    title = "Scatter Plot of Median Household Income vs. 'vaccine' Search by City",
    x = "Median Household Income",
    y = "Search Popularity: vaccine"
  )
```

p6

`geom\_smooth()` using formula = 'y ~ x'



The results from the Pearson correlation test suggest a positive statistically significant correlation between the median household income and both keywords “covid” and “vaccine”. The scatter plot results are consistent with the correlation tests. For “covid”, the majority of the cities with median household income lower than \$10,000 have search popularity centered around 65. They generally have a wider range of search popularity than those with median household income higher than \$10,000. The latter mostly have over 70 searches for “covid”. Based on the plot of median household income and “loans”, as the median household income increases, the search popularity for “loans” seems to increase as well. About half of the cities



with median household income lower than \$125,000 have search popularity below 60, while the majority of those with median household income higher than \$125,000 have search popularity above 60.