

Covid-19-Mask-Usage-Visualization-with-Google-BigQuery-and-Tableau-

After learning a considerable amount of SQL, I began looking for large, realworld databases where I could practice writing more complex queries.

I became acquainted with Google's Data Warehouse BigQuery and it's vast array of public datasets. Given my background in Sociology, I'm no stranger to the American Community Survey ("ACS") the thorough and expansive demographic survey conducted by the United States Census Bureau. The ACS is comprised of numerous variables including: age, race/ethnicity, gender, employment, housing, and many others. It's an excellent source of information for community level demographic characteristics. When evaluated alongside other variables of interest it becomes a treasure trove of explanatory power.

Enter BigQuery.

The ACS collects data at many different levels (state, county, [zipcode](#), geographic coordinates). Since it's county level data uses the county fips code as it's unique primary key, it's easy to combine the ACS table with other tables that have county level data identified with the county fips code.

The New York Times conducted a survey on county level mask data in July:
<https://www.nytimes.com/interactive/2020/07/17/upshot/coronavirus-face-mask-map.html>

The New York Times also regularly collects county level COVID-19 data: <https://www.nytimes.com/interactive/2021/us/covid-cases.html>

Both of these datasets are housed in google's data warehouse and both use the county fips code as their unique primary key.

Using SQL, I wrote this query to retrieve a dataset combining these three data sources.

```

1 SELECT
2   county,
3   state_name,
4   employed_pop/ total_pop AS employed_pop_pct,
5   unemployed_pop/ total_pop AS unemployed_pop_pct,
6   pop_in_labor_force/ total_pop AS pop_in_labor_force_pct,
7   not_in_labor_force/ total_pop AS not_in_labor_force_pct,
8   workers_16_and_over/ total_pop AS workers_16_and_over_pct,
9   armed_forces/ total_pop AS armed_forces_pct,
10  civilian_labor_force/ total_pop AS civilian_labor_force_pct,
11  employed_agriculture_forestry_fishing_hunting_mining/ total_pop AS employed_agriculture_forestry_fishing_hunting_mining_pct,
12  employed_arts_entertainment_recreation_accommodation_food/ total_pop AS employed_arts_entertainment_recreation_accommodation_food_pct,
13  employed_construction/ total_pop AS employed_construction_pct,
14  employed_education_health_social/ total_pop AS employed_education_health_social_pct,
15  employed_finance_insurance_real_estate/ total_pop AS employed_finance_insurance_real_estate_pct,
16  employed_information/ total_pop AS employed_information_pct,
17  employed_manufacturing/ total_pop AS employed_manufacturing_pct,
18  employed_other_services_not_public_admin/ total_pop AS employed_other_services_not_public_admin_pct,
19  employed_public_administration/ total_pop AS employed_public_administration_pct,
20  employed_retail_trade/ total_pop AS employed_retail_trade_pct,
21  employed_science_management_admin_waste/ total_pop AS employed_science_management_admin_waste_pct,
22  employed_transportation_warehousing_utilities/ total_pop AS employed_transportation_warehousing_utilities_pct,
23  employed_wholesale_trade/ total_pop AS employed_wholesale_trade_pct,
24  occupation_management_arts/ total_pop AS occupation_management_arts_pct,
25  occupation_natural_resources_construction_maintenance/ total_pop AS occupation_natural_resources_construction_maintenance_pct,
26  occupation_production_transportation_material/ total_pop AS occupation_production_transportation_material_pct,
27  occupation_sales_office/ total_pop AS occupation_sales_office_pct,
28  occupation_services/ total_pop AS occupation_services_pct,
29  management_business_sci_arts_employed/ total_pop AS management_business_sci_arts_employed_pct,
30  sales_office_employed/ total_pop AS sales_office_employed_pct,
31  poverty/ total_pop AS poverty_pct,
32  gini_index,
33  median_income,
34  median_rent,
35  percent_income_spent_on_rent,
36  million_dollar_housing_units,
37  black_pop/ total_pop AS black_pop_pct,
38  hispanic_pop/ total_pop AS hispanic_pop_pct,
39  asian_pop/ total_pop AS asian_pop_pct,
40  white_pop/ total_pop AS white_pop_pct,
41  other_race_pop/ total_pop AS other_race_pop_pct,
42  amerindian_pop/ total_pop AS amerindian_pop_pct,
43  total_cases / total_pop AS covid_cases_per_capita,
44  total_deaths/ total_pop AS deaths_per_capita,
45  total_pop,
46  mask_usage
47 --Selecting Variables of Interest
48 FROM
49   `bigquery-public-data.census_bureau.acs county_2018_5yr` acs
50 JOIN
51   `bigquery-public-data.census_bureau.acs county_2018_5yr` acs
52 --Joining the ACS data with a subquery that retrieves NYT covid-19 and mask use data
53 (SELECT *,
54  CASE
55  WHEN mask_score < 0.4010546 THEN 'low'
56  WHEN mask_score > 0.4010546 AND mask_score < 0.5788573 THEN 'med'
57  WHEN mask_score > 0.5788573 THEN 'high'
58  END AS mask_usage
59 --Categorizing mask usage by splitting the data into thirds by calculating the probability distribution function in R
60 FROM
61   (SELECT covid_county, covid_state_name, covid_total_cases, covid_total_deaths, covid_county_fips_code,
62    (never * -1) + (rarely * -5) + (sometimes * 5) + (frequently * 5) + (always * 1) AS mask_score
63    --Converting the five question survey data into a single statistic
64    FROM
65     (SELECT covid_county, covid_state_name, covid_fips_code, SUM(confirmed_cases) as total_cases, SUM(deaths) as total_deaths
66      FROM `bigquery-public-data.covid19_nyt_us_counties`
67      GROUP BY covid_county, covid_fips_code) AS covid
68 JOIN `bigquery-public-data.covid19_nyt_mask_use_by_county` AS mask
69 ON covid_county_fips_code = mask_county_fips_code) AS dat1
70 --Joining the NYT Covid-19 data with the mask use data
71 ON
72  dat1.county_fips_code = acs_geo_id

```

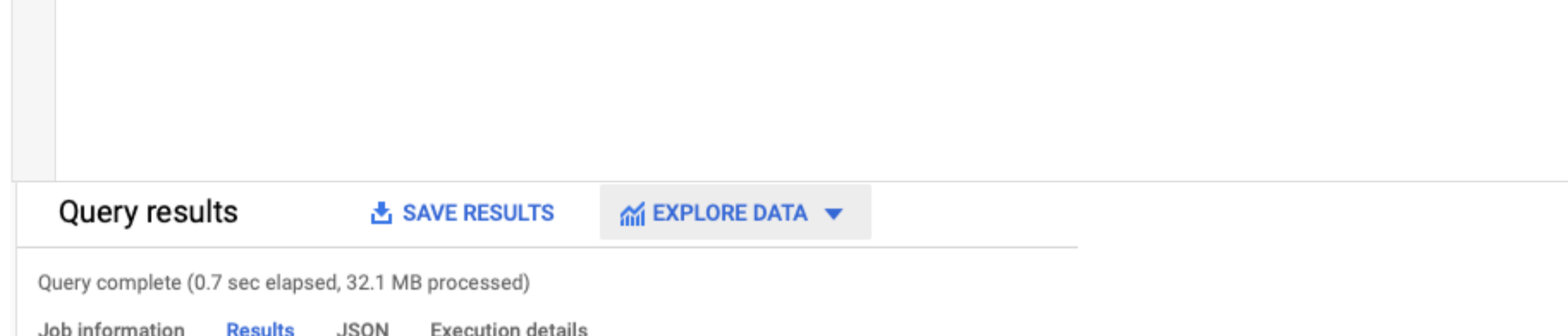
As you can see, I converted the mask usage survey results into a single statistic which I call "mask score".

I also wrote a second query to retrieve summary statistics

```

1 SELECT AVG(mask_score), STDEV_POP(mask_score)
2 FROM
3 (SELECT covid_county, covid_state_name, covid_total_cases, covid_total_deaths, covid_county_fips_code,
4  (never * -1) + (rarely * -5) + (sometimes * 5) + (frequently * 5) + (always * 1) AS mask_score
5 FROM
6  (SELECT covid_county, covid_state_name, covid_fips_code, SUM(confirmed_cases) as total_cases, SUM(deaths) as total_deaths
7   FROM `bigquery-public-data.covid19_nyt_us_counties`
8   GROUP BY covid_county, covid_fips_code) AS covid
9 JOIN `bigquery-public-data.covid19_nyt_mask_use_by_county` AS mask
10 ON covid_county_fips_code = mask_county_fips_code
11 ORDER BY mask_score) AS dat

```



Then I used the qnorm function in r to determine the 33.3 and 66.6 percentile markers.

```

> med_threshold <- qnorm(2/3, mean = 0.48995595276093235, sd = 0.20639825812739154)
>
> med_threshold
[1] 0.5788573
> low_threshold <- qnorm(1/3, mean = 0.48995595276093235, sd = 0.20639825812739154)
>
> low_threshold
[1] 0.4010546

```

I then took the percentile markers and used a Case Statement within my query to classify each county as "high", "med", or "low".

```

54 CASE
55 WHEN mask_score < 0.4010546 THEN 'low'
56 WHEN mask_score > 0.4010546 AND mask_score < 0.5788573 THEN 'med'
57 WHEN mask_score > 0.5788573 THEN 'high'
58 END AS mask_usage

```

I ran my query and returned the following table:

Row	county	state_name	employed_pop_pct	unemployed_pop_pct	pop_in_labor_force_pct	not_in_labor_force_pct	workers_16_and_over_pct	armed_forces_pct	civilian_labor_force_pct	employed_pct
1	Apache	Arizona	0.25677413942563126	0.035024188361623	0.291812304991751	0.4638293112608708	0.25392187019378654	1.398171192080758E-5	0.2917983277872543	

Then I downloaded the data as a csv file and imported it into Tableau for visualization and analysis.

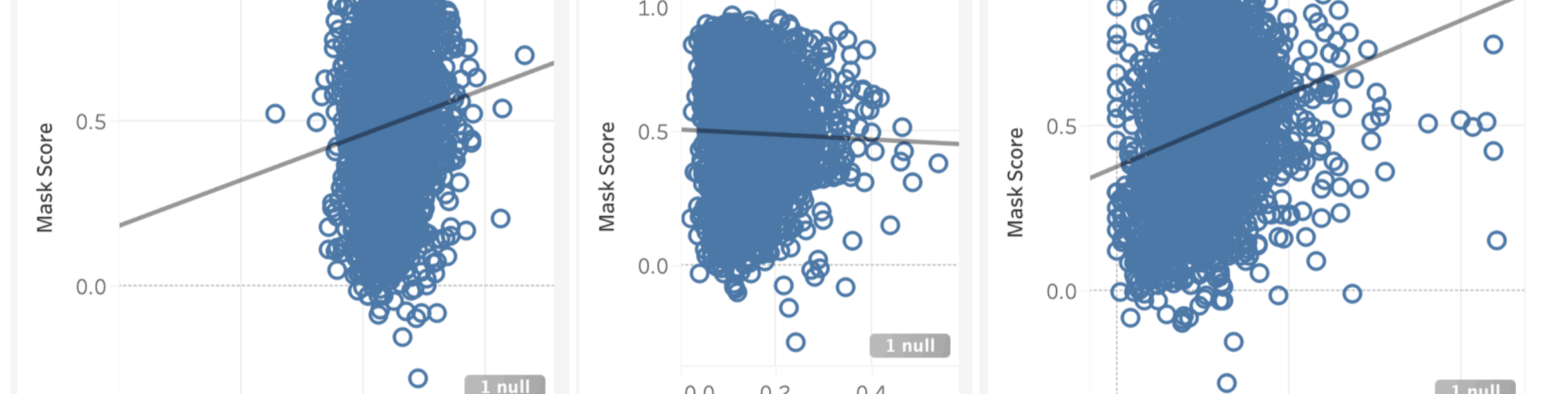
I began by examining the relationships between mask usage and several other variables related to income and employment. I did this by creating a parameter to allow easy adjustment between variables and a calculated field to implement the parameter selection.

```

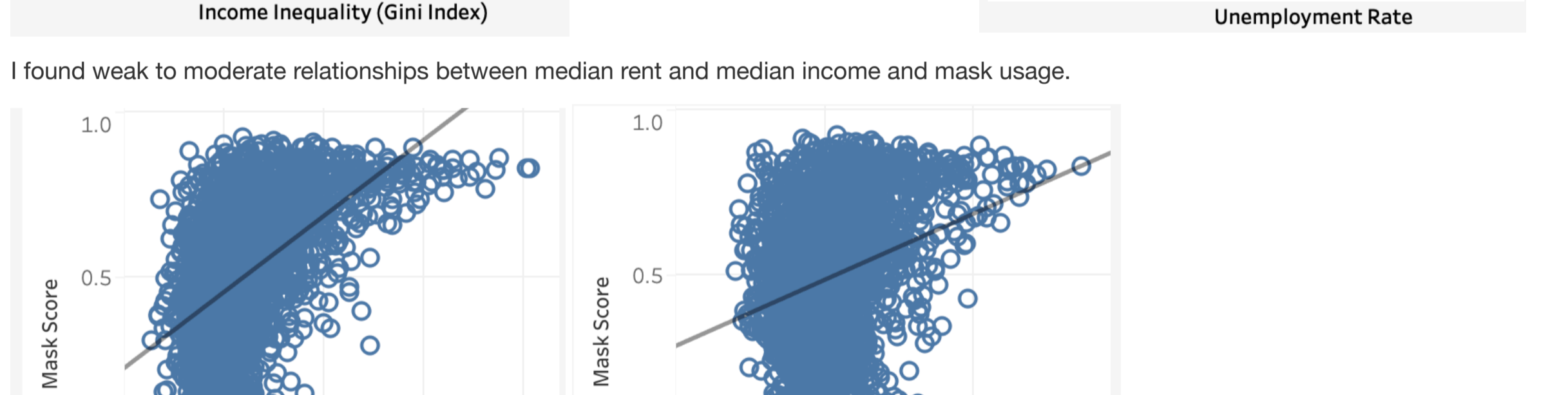
IF [Select Variable] = 1
THEN [Income Inequality]
ELSEIF [Select Variable] = 2
THEN [Median Income]
ELSEIF [Select Variable] = 3
THEN [Median Rent]
ELSEIF [Select Variable] = 4
THEN [Poverty Rate]
ELSEIF [Select Variable] = 5
THEN [Unemployment Rate]
END

```

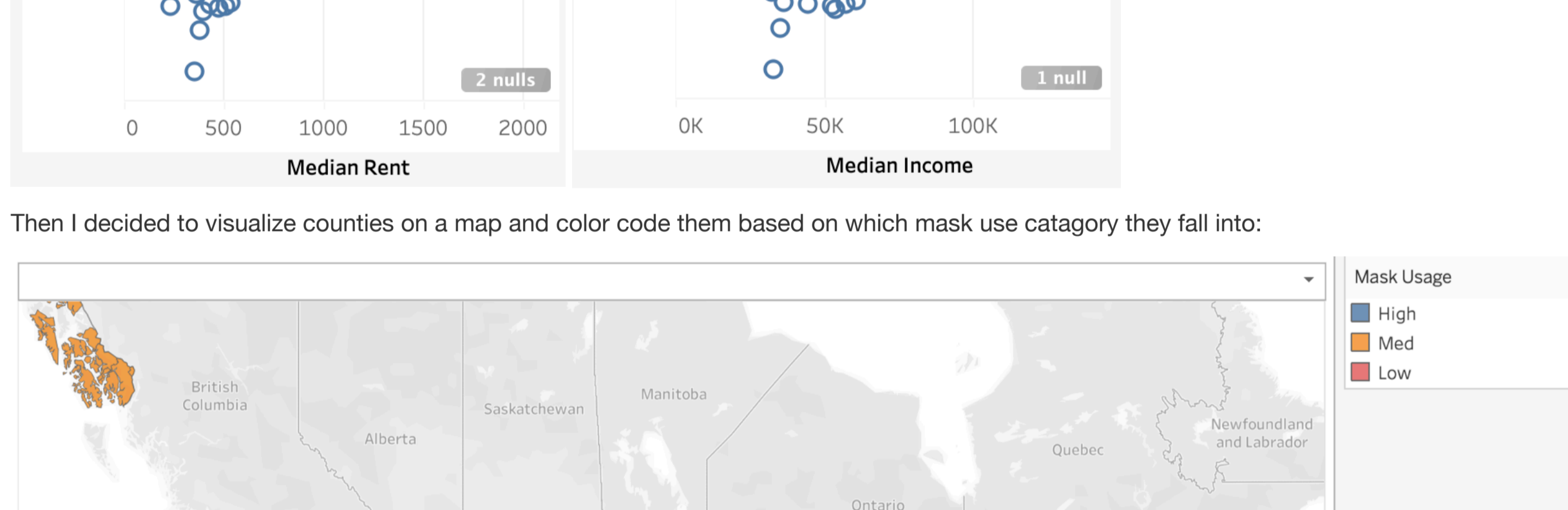
I discovered that there was little to no relationship between income inequality, poverty or unemployment and mask usage.



I found weak to moderate relationships between median rent and median income and mask usage.

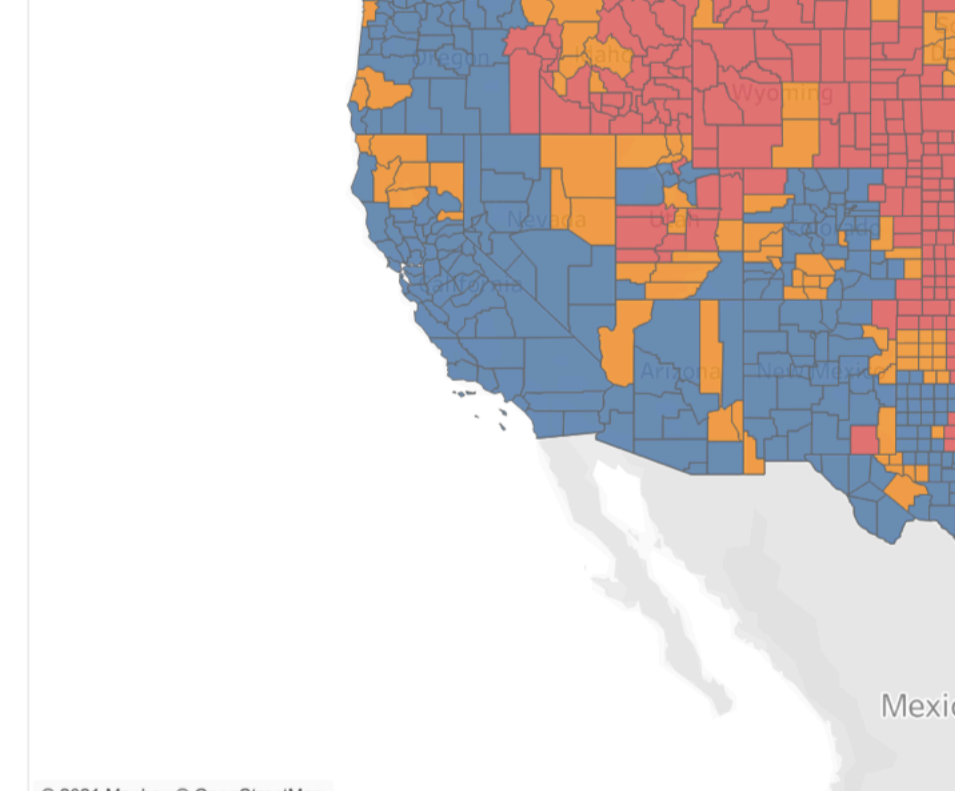


Then I decided to visualize counties on a map and color code them based on which mask use category they fall into:



As you can see mask use is considerably stronger in the coastal areas than in the middle country.

I then turned my attention to diversity. When I saw that the low and medium mask counties had a considerably higher population percentage, I knew I needed to find away to quantify diversity.

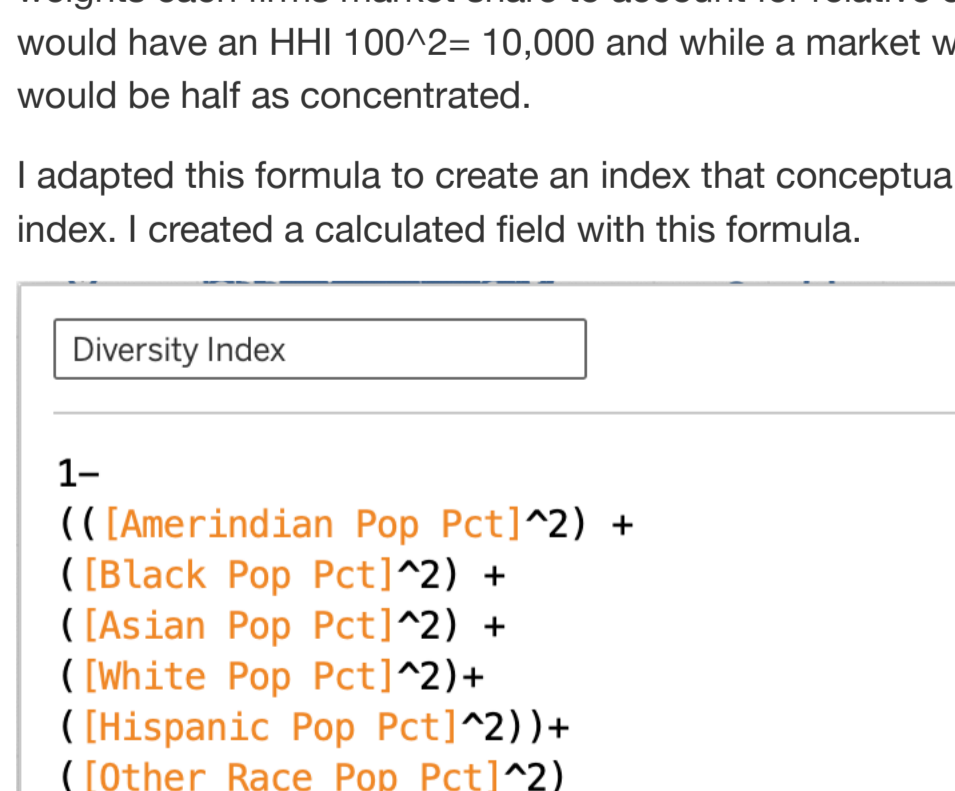


In my day job working at a law firm, I've been tasked with measuring market concentration before and after a potential merger to assess its legality. We use a measure called the Herfindahl-Hirschman Index (HHI). The formula for HHI is $HHI = S_1^2 + S_2^2 + S_3^2 + \dots + S_n^2$. Where S denotes a firm's market share. Each market share is the expressed as the percentage of total market revenue, and squaring each market share weights each firm market share to account for relative concentration. For example a market with one firm controlling 100% of the market share would have an HHI $100^2 = 10,000$ and a market with two firms each controlling 50% would have an HHI of $50^2 + 50^2 = 5,000$ and so would be half as concentrated.

I adapted this formula to create an index that conceptualized how concentrated one racial group was in each county, which I labeled diversity index. I created a calculated field with this formula.

I subtracted the index from one so that higher values are more diverse and not more concentrated in one race since that is what the formula measured.

As you can see, the high mask areas tend to be more diverse and all of the most diverse areas are all high mask.



I conducted exploratory data analysis on numerous other values, but did not find inferences significant enough to warrant space on the dashboard. My next project is to evaluate the relationship between mask wearing and COVID-19.

You can view the completed dashboard here:

<https://public.tableau.com/profile/william.schulman#vizhome/PredictorsofCovid-19MaskWearing/Dashboard1>