

Code - A bit more organized

Dylan Koproski

2025-01-28

Libraries

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2    3.5.1      v tibble    3.2.1
## v lubridate  1.9.3      v tidyr     1.3.1
## v purrr      1.0.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(jsonlite)
```

```
##
## Attaching package: 'jsonlite'
##
## The following object is masked from 'package:purrr':
##
##   flatten
```

```
library(readxl)
library(rsample)
library(caret)
```

```
## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following object is masked from 'package:purrr':
##
##   lift
```

```
library(VGAM)
```

```
## Loading required package: stats4
## Loading required package: splines
##
## Attaching package: 'VGAM'
##
## The following object is masked from 'package:caret':
```

```

##
##   predictors
library(COMPoissonReg)

## Loading required package: Rcpp
##
## Attaching package: 'Rcpp'
##
## The following object is masked from 'package:rsample':
##
##   populate
##
## Loading required package: numDeriv
##
## Attaching package: 'COMPoissonReg'
##
## The following object is masked from 'package:VGAM':
##
##   get.offset
library(pscl)

## Classes and Methods for R originally developed in the
## Political Science Computational Laboratory
## Department of Political Science
## Stanford University (2002-2015),
## by and under the direction of Simon Jackman.
## hurdle and zeroinfl functions by Achim Zeileis.
library(lme4)

## Loading required package: Matrix
##
## Attaching package: 'Matrix'
##
## The following objects are masked from 'package:tidyr':
##
##   expand, pack, unpack
library(zipcodeR)
library(maps)

##
## Attaching package: 'maps'
##
## The following object is masked from 'package:purrr':
##
##   map
library(MASS)

##
## Attaching package: 'MASS'
##
## The following object is masked from 'package:dplyr':
##
##   select

```

```

library(usmap)
library(scales)

##
## Attaching package: 'scales'
##
## The following object is masked from 'package:purrr':
##
##   discard
##
## The following object is masked from 'package:readr':
##
##   col_factor
conflicted::conflict_prefer("select", "dplyr")

## [conflicted] Will prefer dplyr::select over any other package.
conflicted::conflict_prefer("map", "purrr")

## [conflicted] Will prefer purrr::map over any other package.
conflicted::conflict_prefer("filter", "dplyr")

## [conflicted] Will prefer dplyr::filter over any other package.

```

Preprocessing

```

# Load data
df_visitor = read_csv("mobility.csv")

## Warning: One or more parsing issues, call `problems()` on your data frame for details,
## e.g.:
##   dat <- vroom(...)
##   problems(dat)

## Rows: 24583 Columns: 52
## -- Column specification -----
## Delimiter: ","
## chr  (30): placekey, parent_placekey, safegraph_brand_ids, location_name, br...
## dbl  (14): naics_code, latitude, longitude, phone_number, wkt_area_sq_meters...
## lgl   (3): enclosed, is_synthetic, includes_parking_lot
## dtm   (2): date_range_start, date_range_end
## date  (3): opened_on, closed_on, tracking_closed_since
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
df_census = read_csv("tract_census.csv", skip = 1) |>
  janitor::clean_names()

## New names:
## * `` -> `...459`

## Warning: One or more parsing issues, call `problems()` on your data frame for details,
## e.g.:
##   dat <- vroom(...)

```

```

##   problems(dat)

## Rows: 85395 Columns: 459
## -- Column specification -----
## Delimiter: ","
## chr (272): Geography, Geographic Area Name, Estimate!!Total!!Total populatio...
## dbl (186): Estimate!!Total!!Total population, Margin of Error!!Total!!Total ...
## lgl   (1): ...459
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
df_tract_zip = read_excel("tract_zip.xlsx")

# Temporary restriction to NYC
state_county_code_str = c(36005, 36047, 36061, 36081, 36085)

# Process visitor data - none missing
filtered_df_visitor = df_visitor |>
  mutate(first_five_digits = substr(poi_cbg, 1, 5)) |>
  filter(first_five_digits %in% state_county_code_str) |>
  # Dropping missing values for home cbg
  filter(!is.na(visitor_home_aggregation)) |>
  mutate(identifier = row_number()) |>
  mutate(poi_zip = postal_code) |>
  mutate(visitor_home_aggregation = map(visitor_home_aggregation, ~fromJSON(as.character(.)))) |>
  select(location_name, date_range_start, date_range_end, visitor_home_aggregation, top_category, identifier) |>
  unnest_longer(visitor_home_aggregation) |>
  rename(visitor_census_tract = visitor_home_aggregation_id, visitors = visitor_home_aggregation) |>
  mutate(visitors = if_else(visitors == 4, 3, visitors)) |>
  mutate(poi_lat = latitude,
         poi_long = longitude)

# Census data processing
df_census = df_census |>
  rowwise() |>
  mutate(cbg = str_sub(geography, -11)) |>
  #873 locations have an estimated 0 people, we should exclude these.
  filter(estimate_total_total_population > 0)

# Age group proportions in census data, separate into 3 age groups
filtered_df_census_totals =
  df_census |>
  rowwise() |>
  select(estimate_total_total_population, cbg, geographic_area_name, starts_with("estimate")) |>
  mutate(
    total_under_18 = sum(estimate_total_total_population_age_under_5_years,
                        estimate_total_total_population_age_5_to_9_years,
                        estimate_total_total_population_age_10_to_14_years,
                        estimate_total_total_population_age_15_to_19_years) / estimate_total_total_population,
    total_19_65 = sum(estimate_total_total_population_age_20_to_24_years,
                     estimate_total_total_population_age_25_to_29_years,

```

```

        estimate_total_total_population_age_30_to_34_years,
        estimate_total_total_population_age_35_to_39_years,
        estimate_total_total_population_age_40_to_44_years,
        estimate_total_total_population_age_45_to_49_years,
        estimate_total_total_population_age_50_to_54_years,
        estimate_total_total_population_age_55_to_59_years,
        estimate_total_total_population_age_60_to_64_years,
        estimate_total_total_population_age_65_to_69_years) / estimate_total_total_popula

    total_65_plus = sum(estimate_total_total_population_age_70_to_74_years,
        estimate_total_total_population_age_75_to_79_years,
        estimate_total_total_population_age_80_to_84_years,
        estimate_total_total_population_age_85_years_and_over) / estimate_total_total_p

) |>
rename("total" = estimate_total_total_population) |>
select(cbg, geographic_area_name, total, total_under_18, total_19_65, total_65_plus)

# Define primary ZIP code per census tract
primary_tract_zip = df_tract_zip |>
  group_by(tract) |>
  summarize(zip = min(zip)) # Selects the minimum ZIP as primary for simplicity

# Merge filtered_df_census_totals with filtered_df_visitor
merged_df = filtered_df_visitor |>
  inner_join(filtered_df_census_totals, by = c("visitor_census_tract" = "cbg")) |>
  mutate(
    visitors_under_18 = visitors * total_under_18,
    visitors_19_65 = visitors * total_19_65,
    visitors_65_plus = visitors * total_65_plus
  )

# Map primary ZIP codes by merging with primary_tract_zip on the census tract
final_df = merged_df |>
  left_join(primary_tract_zip, by = c("visitor_census_tract" = "tract")) |>
  select(location_name, date_range_start, date_range_end, top_category,
    identifier, poi_cbg, poi_zip, visitors, visitors_under_18,
    visitors_19_65, visitors_65_plus, zip, poi_long, poi_lat) |>
  mutate(visitor_zip = zip)

vis_zip_lat_long = geocode_zip(final_df$visitor_zip)

final_df = final_df |>
  left_join(vis_zip_lat_long, by = join_by(visitor_zip == zipcode)) |>
  mutate(vis_lat = lat,
    vis_long = lng) |>
  select(-lat, -lng)

# Rounding, can be adjusted at will
final_df = final_df |>
  mutate(visitors_under_18 = ceiling(visitors_under_18),
    visitors_19_65 = ceiling(visitors_19_65),
    visitors_65_plus = ceiling(visitors_65_plus)) |>
  mutate(total_visitors = visitors_under_18 + visitors_19_65 + visitors_65_plus)

```

```

# There are 2 rows in the above data that have missing values, the visitor zip is missing. We should fi.

# Long dataframe for modelling purposes
df_long = final_df |>
  pivot_longer(
    cols = starts_with("visitors_"),
    names_to = "age_group",
    names_prefix = "visitors_",
    values_to = "visitor_count"
  ) |>
  mutate(top_category = factor(top_category),
         age_group = factor(age_group, levels = c("under_18", "19_65", "65_plus")))

# Data quality looks good, 2 missing zip codes may need to be handled, but data is large enough maybe d

```

Category definitions

```

all_cat = c(
  "Accounting, Tax Preparation, Bookkeeping, and Payroll Services",
  "Activities Related to Credit Intermediation",
  "Activities Related to Real Estate",
  "Advertising, Public Relations, and Related Services",
  "Agencies, Brokerages, and Other Insurance Related Activities",
  "Architectural, Engineering, and Related Services",
  "Automotive Parts, Accessories, and Tire Stores",
  "Automotive Repair and Maintenance",
  "Bakeries and Tortilla Manufacturing",
  "Beer, Wine, and Liquor Stores",
  "Book Stores and News Dealers",
  "Building Equipment Contractors",
  "Building Finishing Contractors",
  "Building Material and Supplies Dealers",
  "Child Day Care Services",
  "Clothing Stores",
  "Consumer Goods Rental",
  "Couriers and Express Delivery Services",
  "Depository Credit Intermediation",
  "Drinking Places (Alcoholic Beverages)",
  "Drycleaning and Laundry Services",
  "Electronic and Precision Equipment Repair and Maintenance",
  "Electronics and Appliance Stores",
  "Elementary and Secondary Schools",
  "Florists",
  "Furniture Stores",
  "Gasoline Stations",
  "General Medical and Surgical Hospitals",
  "General Merchandise Stores, including Warehouse Clubs and Supercenters",
  "Glass and Glass Product Manufacturing",
  "Grocery Stores",
  "Health and Personal Care Stores",
  "Home Furnishings Stores",
  "Investigation and Security Services",

```

```

"Jewelry, Luggage, and Leather Goods Stores",
"Justice, Public Order, and Safety Activities",
"Legal Services",
"Machinery, Equipment, and Supplies Merchant Wholesalers",
"Museums, Historical Sites, and Similar Institutions",
"Offices of Dentists",
"Offices of Other Health Practitioners",
"Offices of Physicians",
"Offices of Real Estate Agents and Brokers",
"Other Amusement and Recreation Industries",
"Other Financial Investment Activities",
"Other Miscellaneous Manufacturing",
"Other Miscellaneous Store Retailers",
"Other Personal Services",
"Other Professional, Scientific, and Technical Services",
"Other Schools and Instruction",
"Other Specialty Trade Contractors",
"Personal and Household Goods Repair and Maintenance",
"Personal Care Services",
"Printing and Related Support Activities",
"Promoters of Performing Arts, Sports, and Similar Events",
"Radio and Television Broadcasting",
"Religious Organizations",
"Restaurants and Other Eating Places",
"Shoe Stores",
"Sound Recording Industries",
"Special Food Services",
"Specialized Design Services",
"Specialty (except Psychiatric and Substance Abuse) Hospitals",
"Specialty Food Stores",
"Sporting Goods, Hobby, and Musical Instrument Stores",
"Support Activities for Road Transportation",
"Technical and Trade Schools",
"Transit and Ground Passenger Transportation",
"Traveler Accommodation",
"Warehousing and Storage",
"Wired and Wireless Telecommunications Carriers"
)

medical_services = c(
  "General Medical and Surgical Hospitals",
  "Health and Personal Care Stores",
  "Offices of Dentists",
  "Offices of Other Health Practitioners",
  "Specialty (except Psychiatric and Substance Abuse) Hospitals",
  "Offices of Physicians"
)

essential_services = c(
  "Health and Personal Care Stores",
  "Pharmacies and Drug Stores",
  "Grocery Stores",
  "Gasoline Stations",

```

```

    "Depository Credit Intermediation",
    "Public Transport Hubs",
    "Government Offices"
)

retail_shopping = c(
    "General Merchandise Stores, including Warehouse Clubs and Supercenters",
    "Clothing Stores",
    "Shoe Stores",
    "Jewelry, Luggage, and Leather Goods Stores",
    "Electronics and Appliance Stores",
    "Furniture Stores",
    "Home Furnishings Stores",
    "Specialty Food Stores",
    "Sporting Goods, Hobby, and Musical Instrument Stores",
    "Book Stores and News Dealers"
)

entertainment_recreation = c(
    "Other Amusement and Recreation Industries",
    "Museums, Historical Sites, and Similar Institutions",
    "Promoters of Performing Arts, Sports, and Similar Events",
    "Radio and Television Broadcasting",
    "Sound Recording Industries"
)

personal_services = c(
    "Personal Care Services",
    "Drycleaning and Laundry Services",
    "Other Personal Services",
    "Personal and Household Goods Repair and Maintenance"
)

hospitality_lodging = c(
    "Traveler Accommodation",
    "Bed and Breakfast Inns",
    "Resorts",
    "Extended Stay Hotels"
)

office_professional = c(
    "Accounting, Tax Preparation, Bookkeeping, and Payroll Services",
    "Legal Services",
    "Architectural, Engineering, and Related Services",
    "Agencies, Brokerages, and Other Insurance Related Activities",
    "Offices of Physicians",
    "Offices of Dentists",
    "Offices of Other Health Practitioners",
    "Real Estate Agencies"
)

target_categories = c("Drinking Places (Alcoholic Beverages)", "Restaurants and Other Eating Places", "I

```



```
df_long_model_filtered_1 = df_long |>
  mutate(non_restaurant = if_else(top_category %in% target_categories, "No", "Yes"))
```

Exploratory Data Analysis (EDA)

Zip Code Flow Matrix

```
# Remove rows with NA in visitor_zip or poi_zip
zip_matrix = final_df |>
  filter(!is.na(visitor_zip) & !is.na(poi_zip)) |>
  group_by(visitor_zip, poi_zip) |>
  summarize(total_visitors = sum(visitors, na.rm = TRUE)) |>
  pivot_wider(names_from = poi_zip, values_from = total_visitors, values_fill = 0)
```

`summarise()` has grouped output by 'visitor_zip'. You can override using the
`.groups` argument.

```
# Convert to matrix and plot
zip_matrix_plot = zip_matrix |>
  column_to_rownames("visitor_zip") |>
  as.matrix() |>
  heatmap(
    col = colorRampPalette(c("white", "red"))(100),
    scale = "none",
    main = "Zip-to-Zip Visitor Flow",
    xlab = "Destination ZIP (To)",
    ylab = "Origin ZIP (From)"
  )
```

State to Destination ZIP in NYC

```
final_df_with_state = final_df |>
  left_join(df_tract_zip, by = c("visitor_zip" = "zip")) |>
  select(!zip) |>
  rename(visitor_state = usps_zip_pref_state)
```

Warning in left_join(final_df, df_tract_zip, by = c(visitor_zip = "zip")): Detected an unexpected many-to-many relationship between 'visitor_zip' and 'zip'.
i Row 1 of `x` matches multiple rows in `y`.
i Row 103268 of `y` matches multiple rows in `x`.
i If a many-to-many relationship is expected, set `relationship = "many-to-many"` to silence this warning.

```
state_zip_matrix = final_df_with_state |>
  group_by(visitor_state, poi_zip) |>
  summarize(total_visitors = sum(visitors, na.rm = TRUE))
```

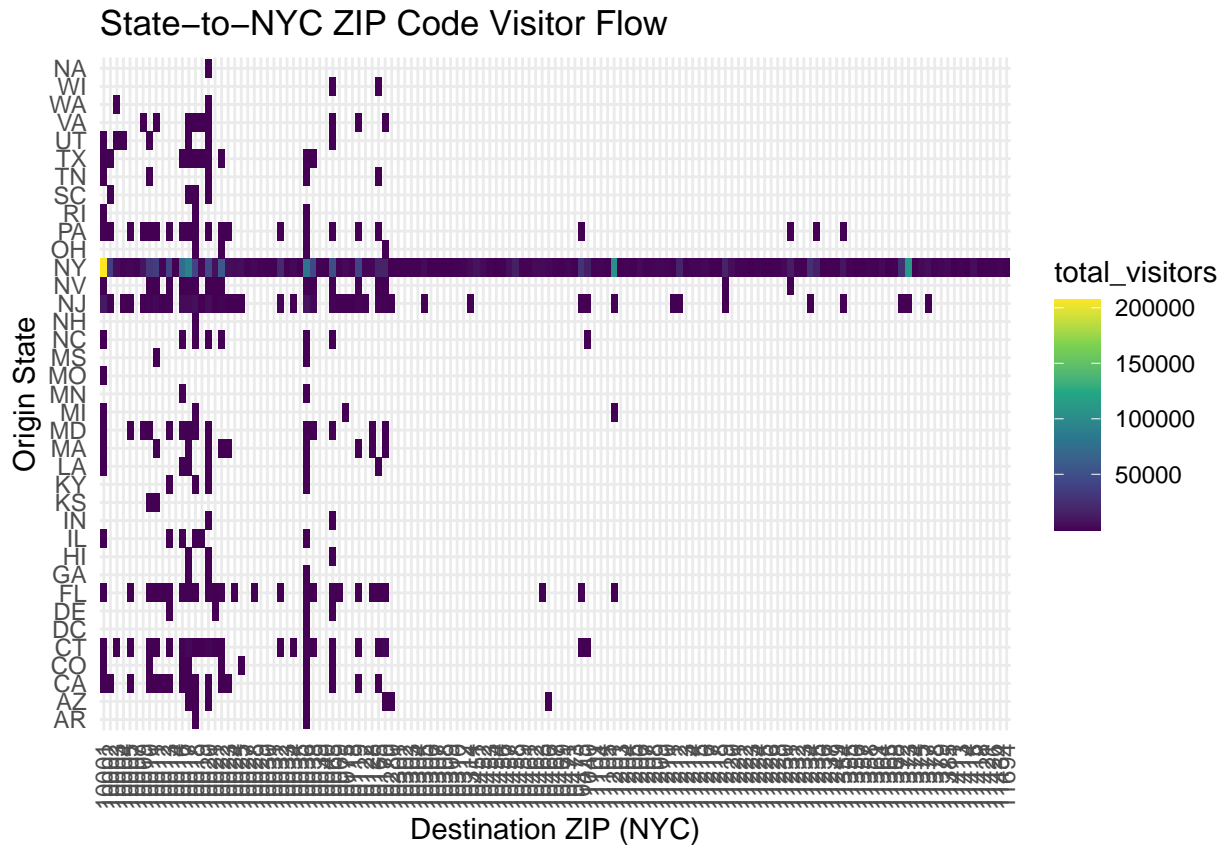
`summarise()` has grouped output by 'visitor_state'. You can override using the
`.groups` argument.

```
ggplot(state_zip_matrix, aes(x = poi_zip, y = visitor_state, fill = total_visitors)) +
  geom_tile() +
  scale_fill_viridis_c() +
  labs(
    title = "State-to-NYC ZIP Code Visitor Flow",
```

```

x = "Destination ZIP (NYC)",
y = "Origin State"
) +
theme_minimal() +
theme(axis.text.x = element_text(angle = 90, vjust = 0.5))

```



```

nyc_zip_visitors = final_df |>
  group_by(poi_long, poi_lat) |>
  summarize(total_visitors = sum(visitors, na.rm = TRUE))

ggplot(nyc_zip_visitors, aes(x = poi_long, y = poi_lat)) +
  stat_density_2d(aes(fill = after_stat(level)), geom = "polygon", color = NA) +
  scale_fill_viridis_c() +
  labs(
    title = "Heatmap of NYC POIs by Visitor Counts",
    x = "Longitude",
    y = "Latitude"
  ) +
  theme_minimal()

```

Flow map - work in progress

```

#excise visitors from hawaii, include the second filter argument to get a better look at lower ends of
flow_df = final_df |>
  filter(vis_long > -150) #/> filter(total_visitors < 15)

```

```

# if you want to include visitors from hawaii
# flow_df = final_df

usa = map_data("state")

usa = rename(usa, state = "region")

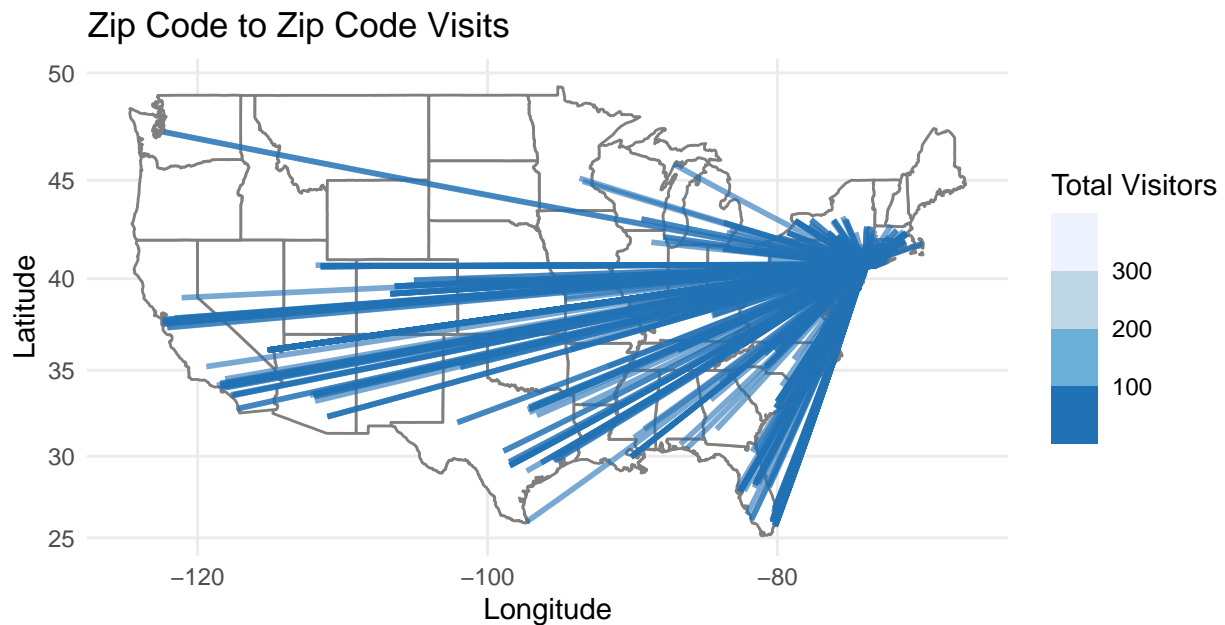
usa$state = str_to_title(usa$state)

stateData = usa |>
  arrange(state, group, order)

ggplot() +
  geom_polygon(data = stateData,
    aes(x = long, y = lat, group = group),
    fill = "white", color = "gray50") +

  geom_segment(data = flow_df,
    aes(x = vis_long, y = vis_lat,
      xend = poi_long, yend = poi_lat,
      color = total_visitors),
    alpha = 0.6, linewidth = 1) +
  scale_color_fermenter(name = "Total Visitors", direction = -1) +
  coord_map() +
  theme_minimal() +
  labs(color = "Volume of Visits",
    title = "Zip Code to Zip Code Visits",
    x = "Longitude",
    y = "Latitude")

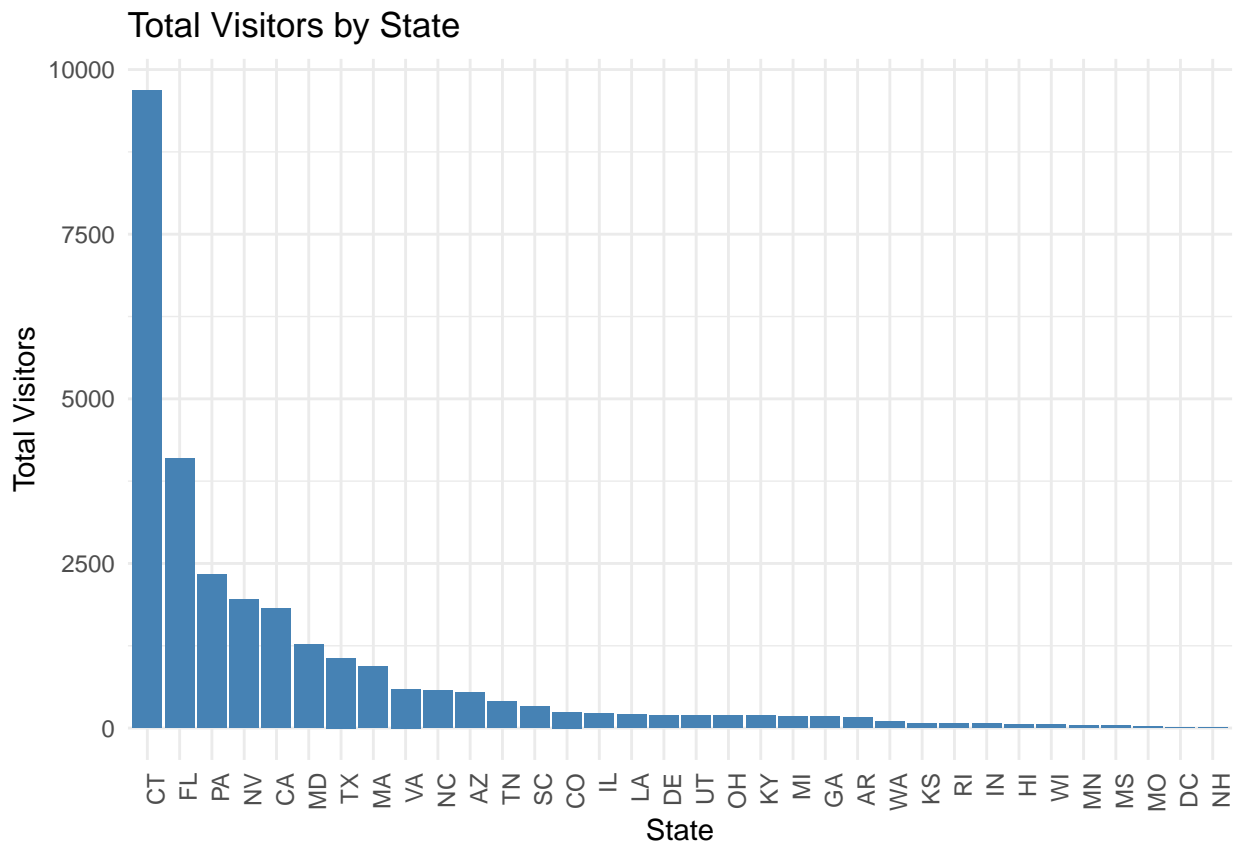
```



Visitor counts

```
state_visitors = final_df_with_state |>
  filter(visitor_state != "NY") |>
  filter(visitor_state != "NJ") |>
  group_by(visitor_state) |>
  summarize(total_visitors = sum(visitors, na.rm = TRUE)) |>
  arrange(desc(total_visitors))

ggplot(state_visitors, aes(x = reorder(visitor_state, -total_visitors), y = total_visitors)) +
  geom_col(fill = "steelblue") +
  labs(
    title = "Total Visitors by State",
    x = "State",
    y = "Total Visitors"
  ) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 90, hjust = 1))
```



Map view of visitor counts

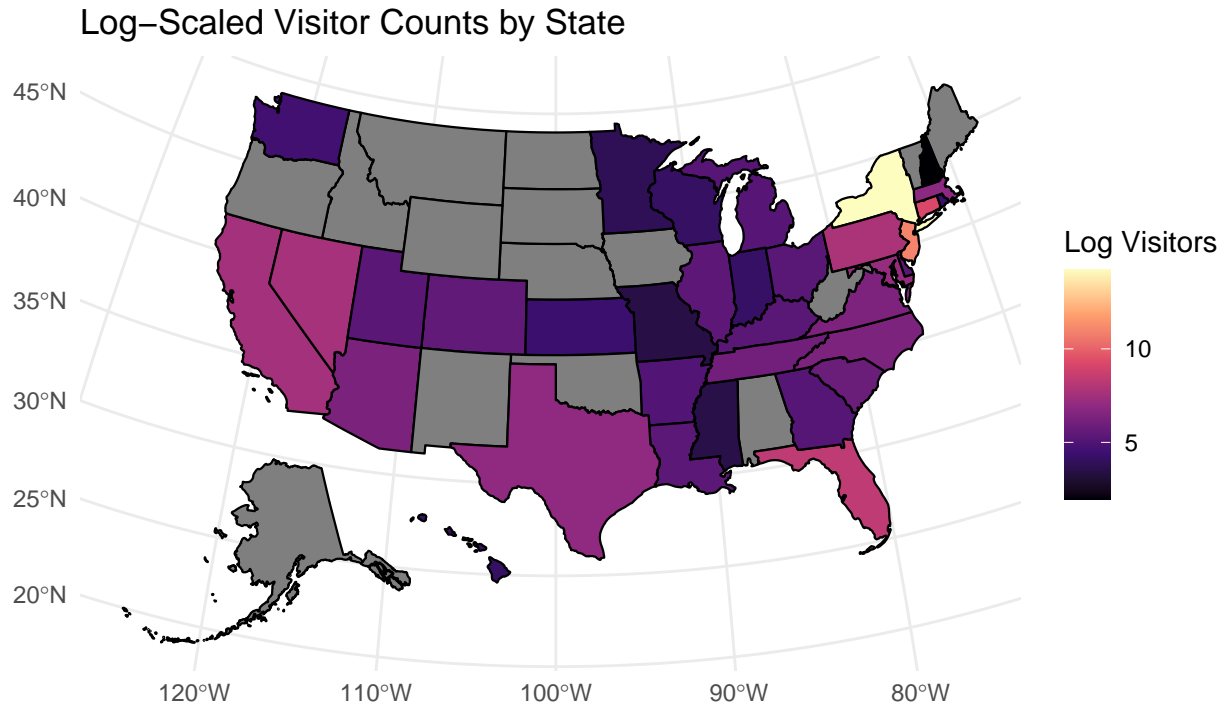
```
# Aggregate visitor counts by state
state_visitors_map = final_df_with_state |>
  group_by(visitor_state) |>
  summarize(total_visitors = sum(visitors, na.rm = TRUE)) |>
  filter(!is.na(visitor_state)) |>
  mutate(log_visitors = log1p(total_visitors)) |>
```

```

rename(state = visitor_state)

# Plot using log scale
plot_usmap(data = state_visitors_map, regions = "states", values = "log_visitors") +
  scale_fill_viridis_c(name = "Log Visitors", option = "magma") +
  labs(title = "Log-Scaled Visitor Counts by State") +
  theme_minimal()

```



```

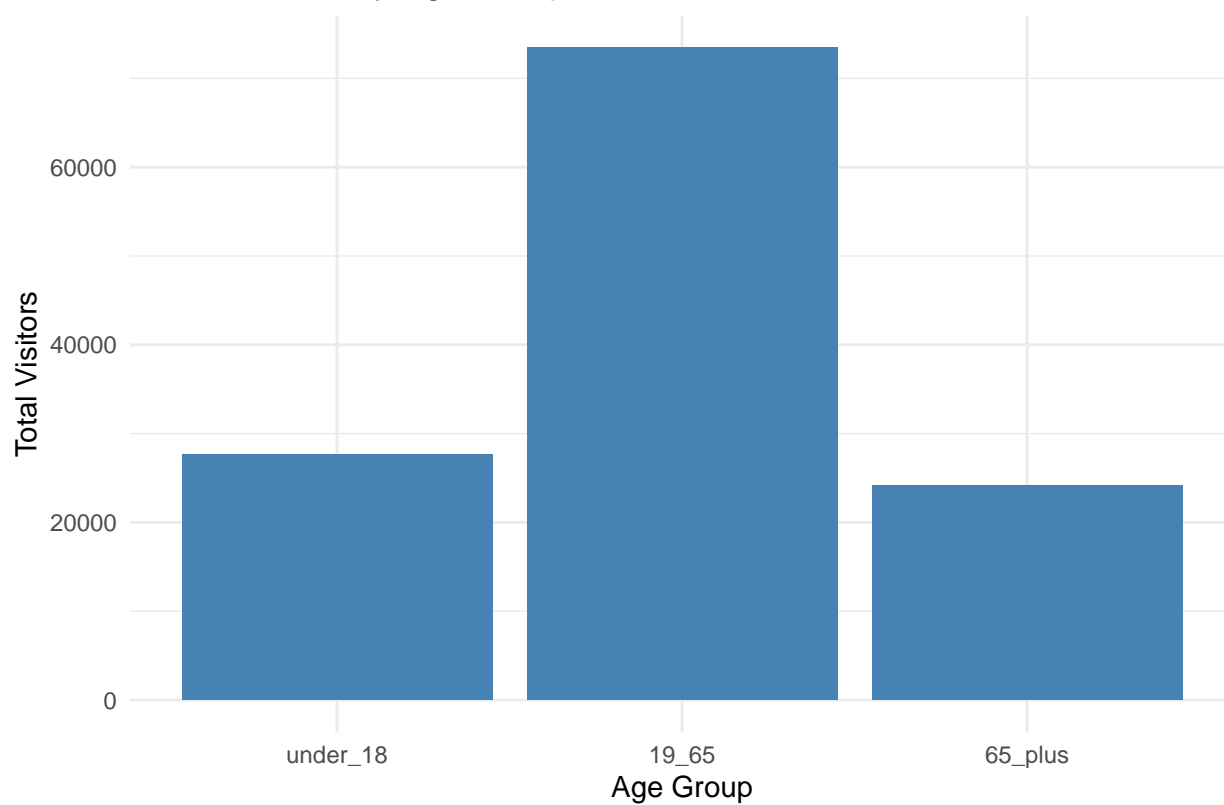
# Function to create bar plots for each category
plot_age_group_counts = function(df, category_name, category_vector) {
  df_filtered = df |>
    filter(top_category %in% category_vector) |>
    group_by(age_group) |>
    summarize(total_visitors = sum(visitor_count, na.rm = TRUE), .groups = "drop")

  ggplot(df_filtered, aes(x = age_group, y = total_visitors)) +
    geom_col(fill = "steelblue") +
    labs(
      title = paste("Visitor Counts by Age Group -", category_name),
      x = "Age Group",
      y = "Total Visitors"
    ) +
    theme_minimal() +
    theme(
      axis.text.x = element_text(angle = 0, hjust = 0.5),
      legend.position = "none" # Remove legend
    )
}

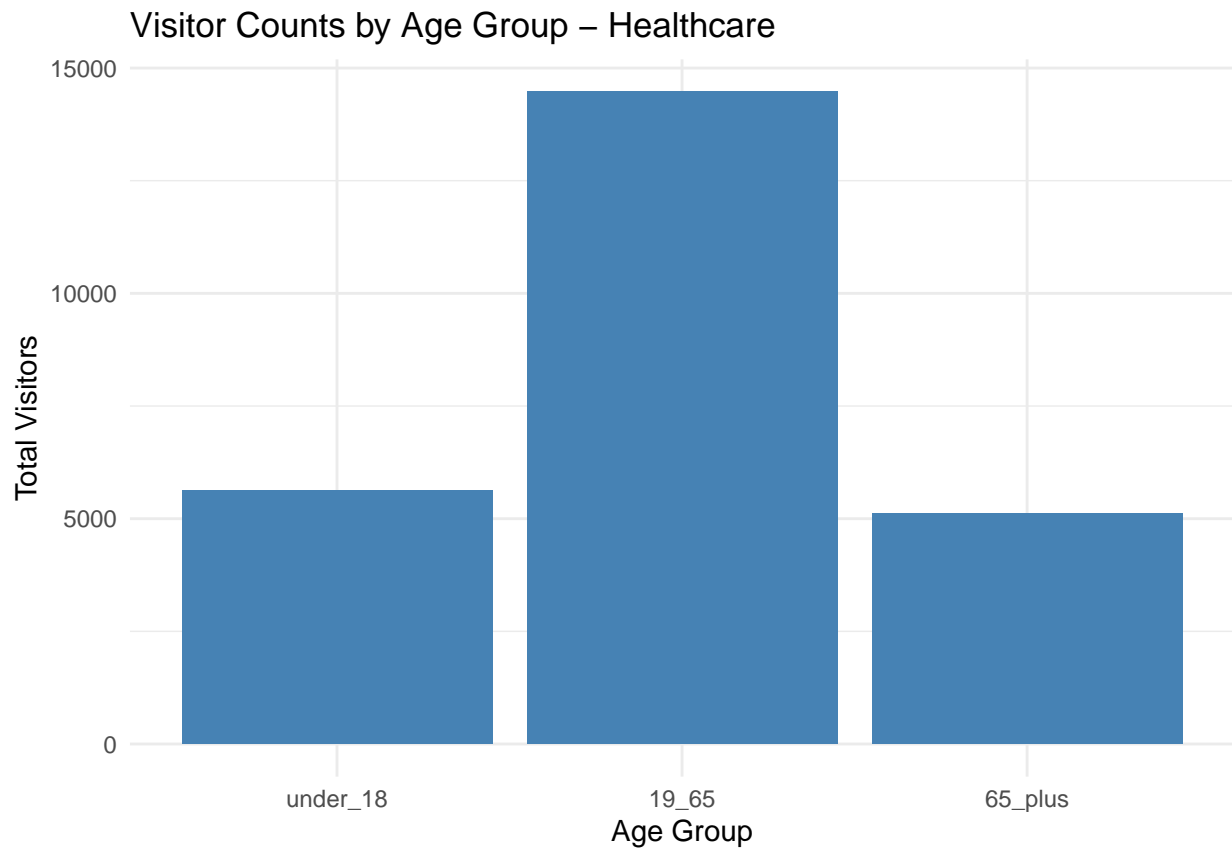
# Generate bar plots for each category
plot_age_group_counts(df_long, "All", all_cat)

```

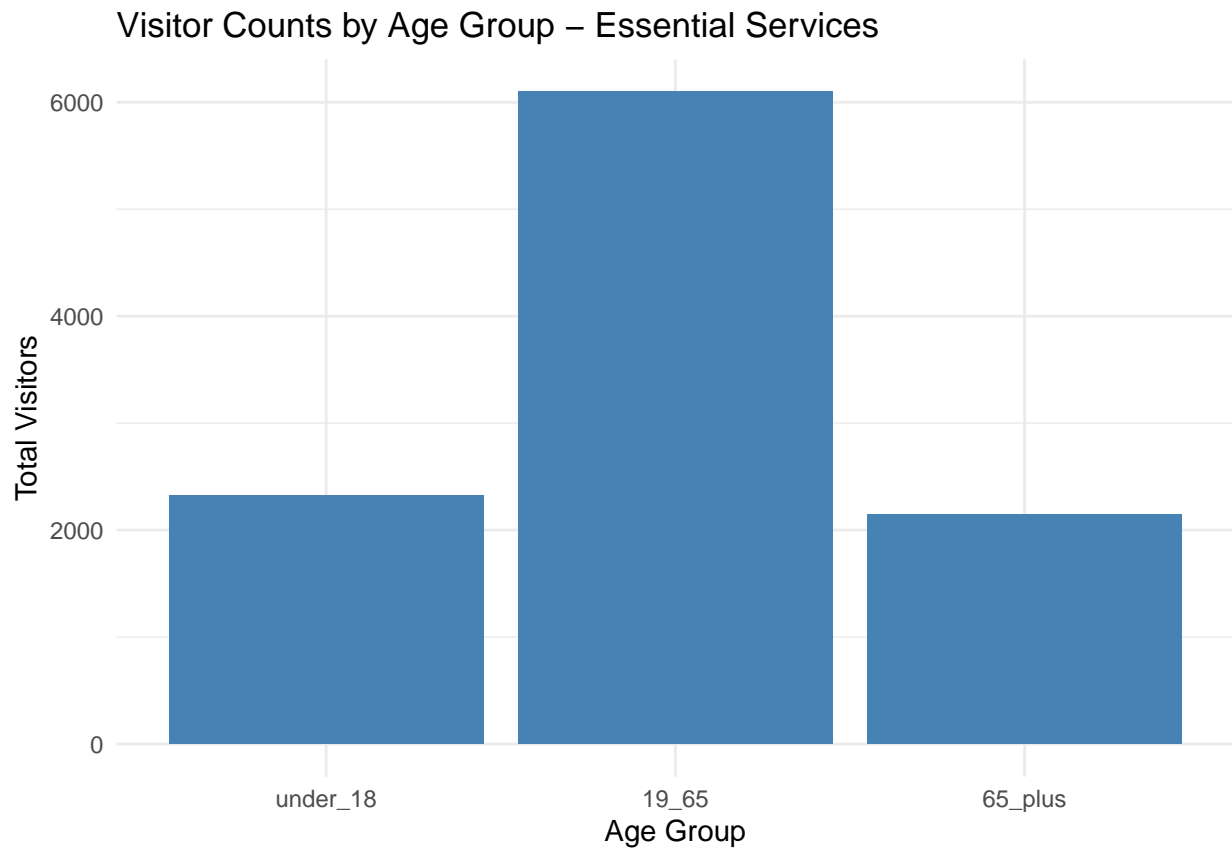
Visitor Counts by Age Group – All



```
plot_age_group_counts(df_long, "Healthcare", medical_services)
```

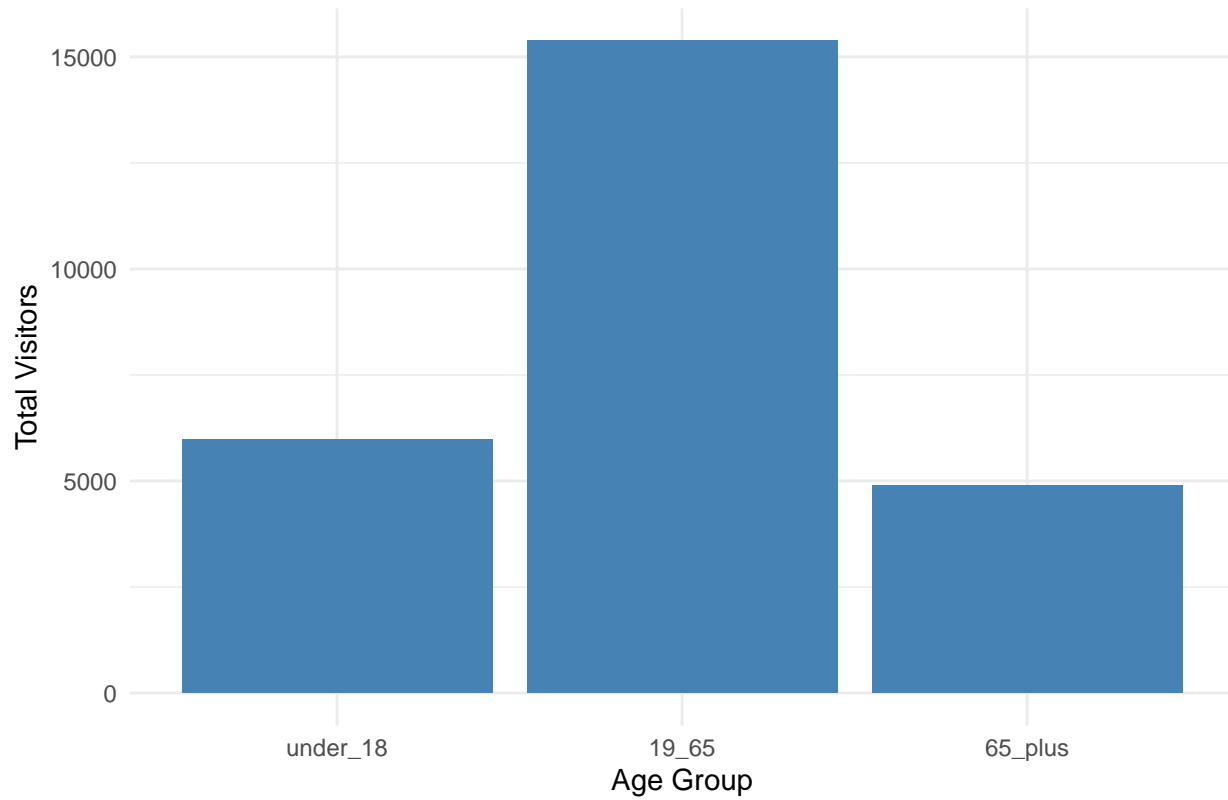


```
plot_age_group_counts(df_long, "Essential Services", essential_services)
```

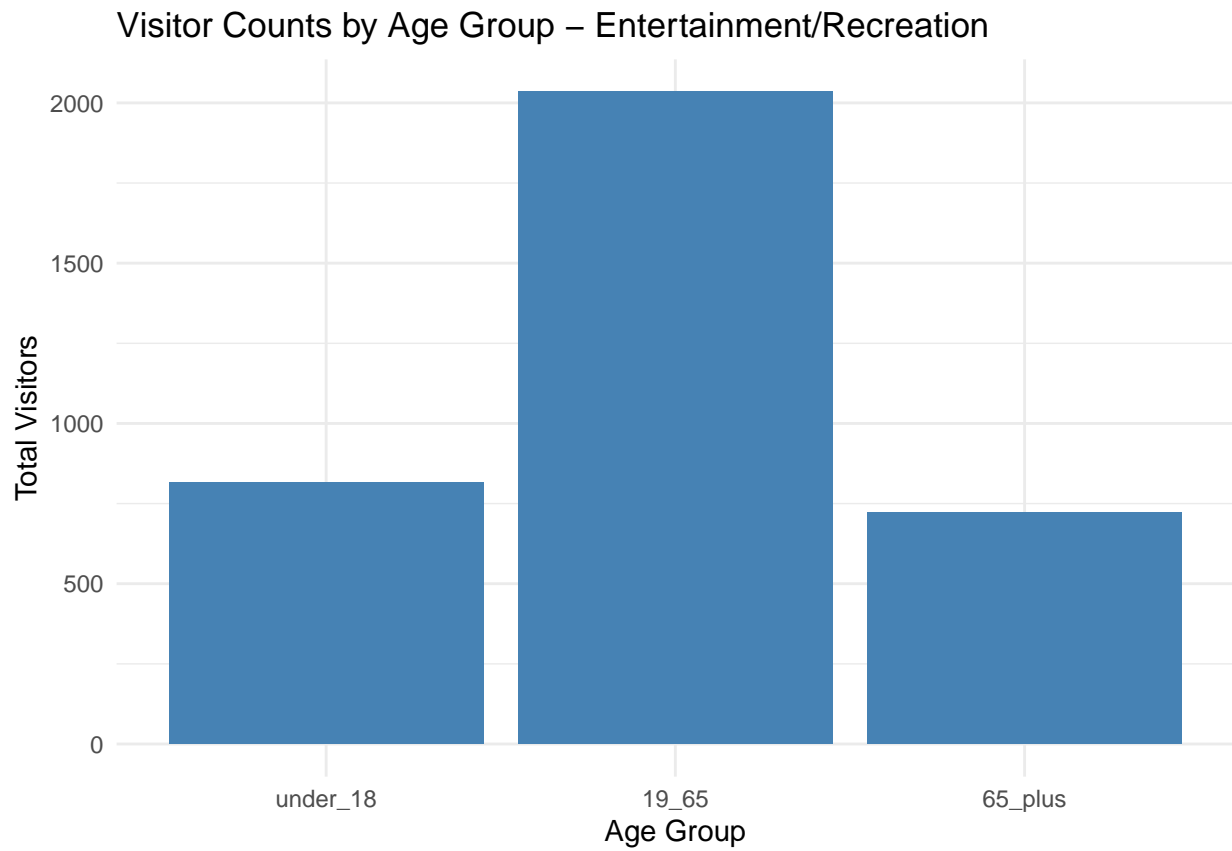


```
plot_age_group_counts(df_long, "Retail Shopping", retail_shopping)
```


Visitor Counts by Age Group – Retail Shopping

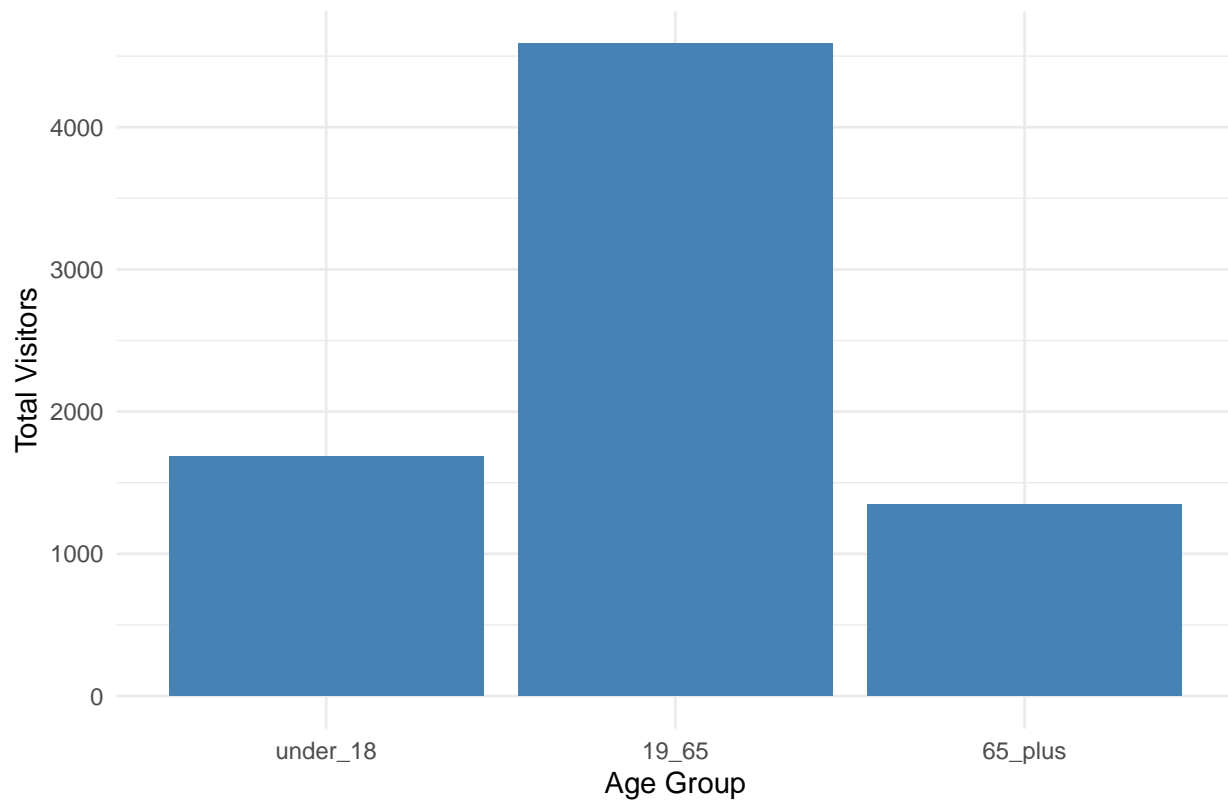


```
plot_age_group_counts(df_long, "Entertainment/Recreation", entertainment_recreation)
```



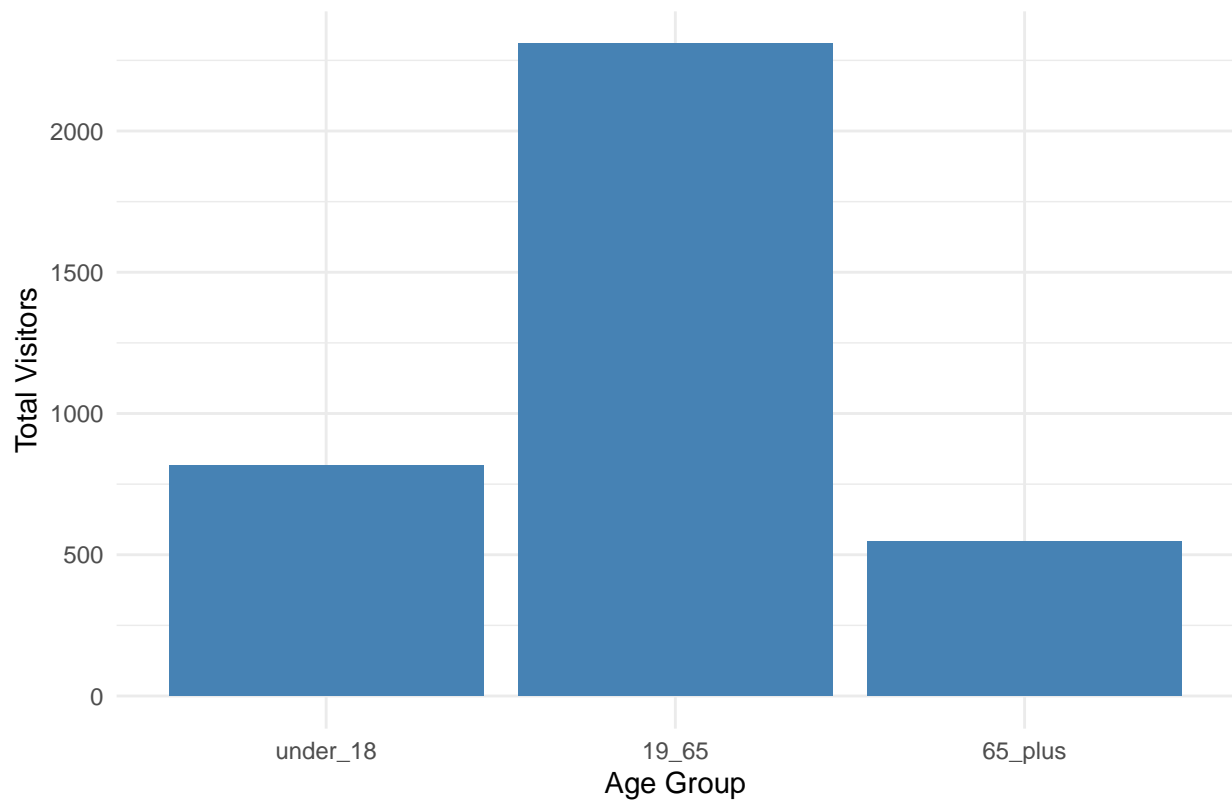
```
plot_age_group_counts(df_long, "Personal Services", personal_services)
```

Visitor Counts by Age Group – Personal Services



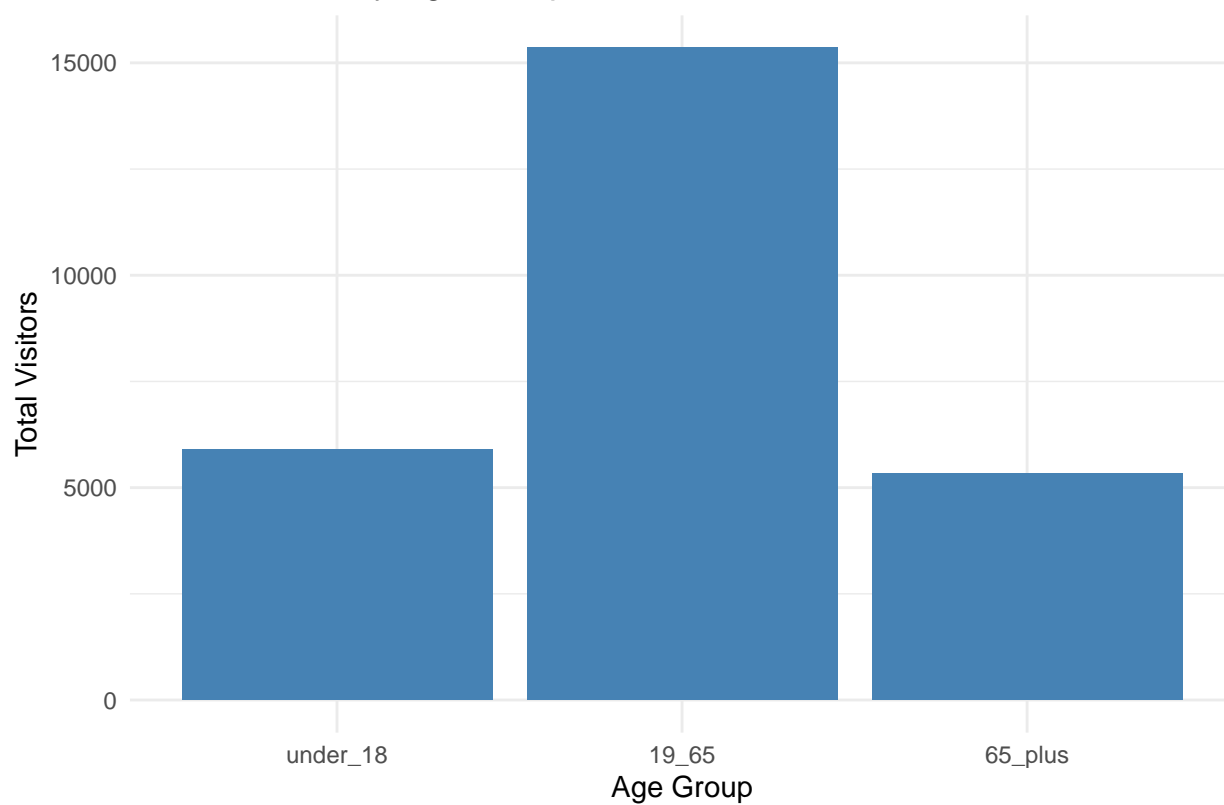
```
plot_age_group_counts(df_long, "Hospitality/Lodging", hospitality_lodging)
```

Visitor Counts by Age Group – Hospitality/Lodging



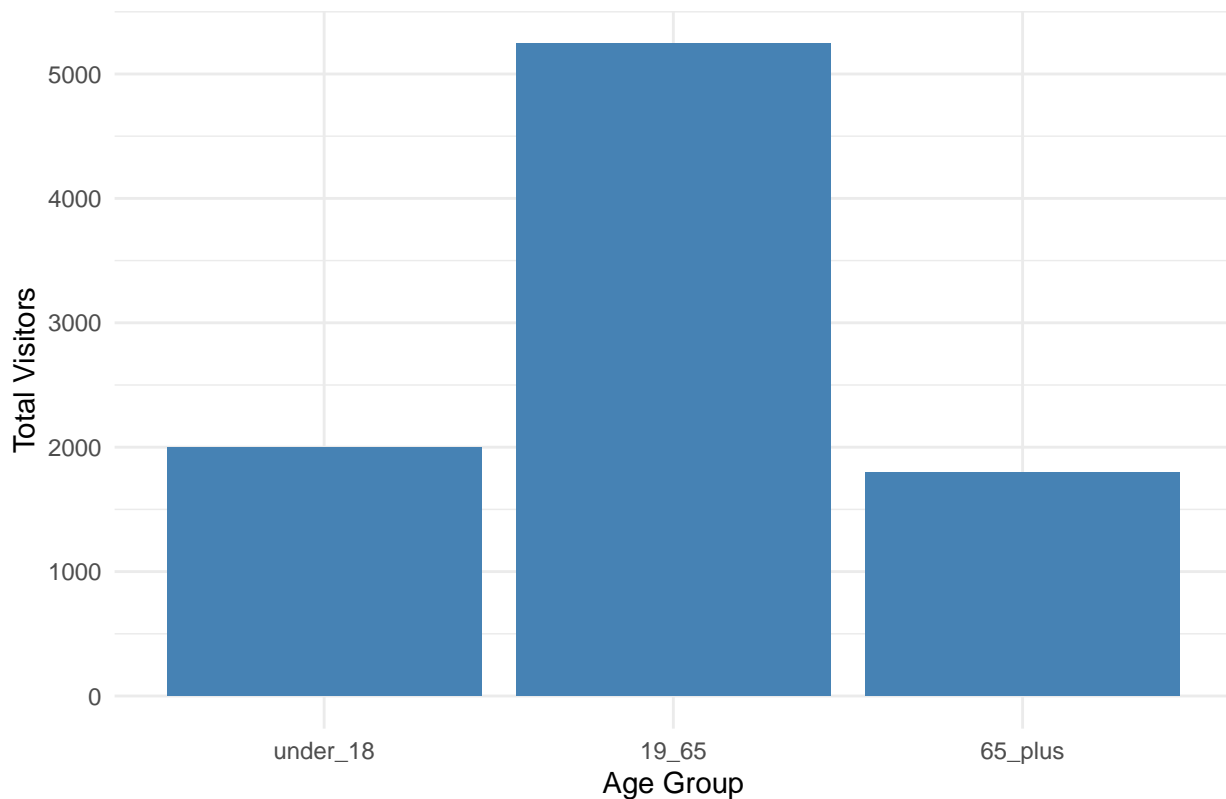
```
plot_age_group_counts(df_long, "Office/Professional", office_professional)
```

Visitor Counts by Age Group – Office/Professional



```
plot_age_group_counts(df_long, "Restaurant/Bar", target_categories)
```

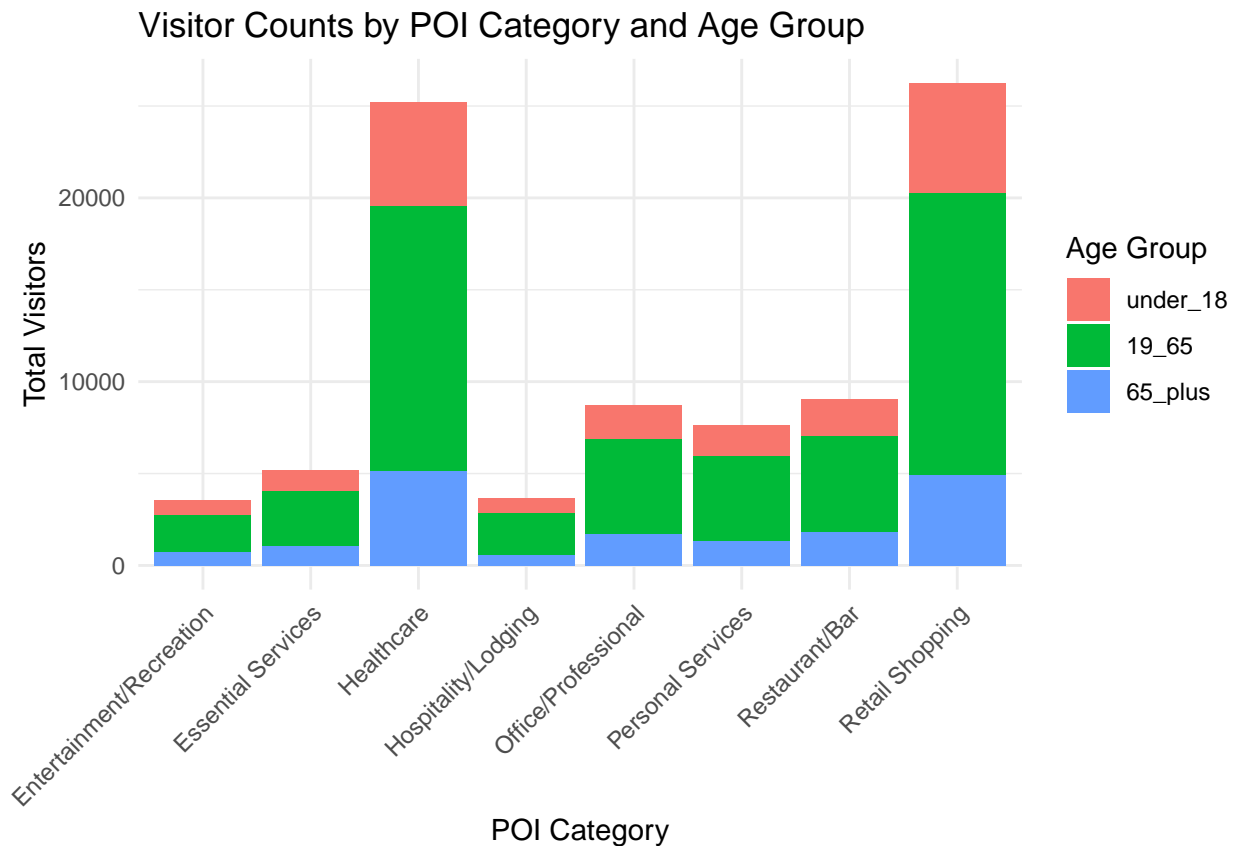
Visitor Counts by Age Group – Restaurant/Bar



```
# Aggregate visitor counts by age group and POI category
df_filtered = df_long |>
  mutate(category_group = case_when(
    top_category %in% medical_services ~ "Healthcare",
    top_category %in% essential_services ~ "Essential Services",
    top_category %in% retail_shopping ~ "Retail Shopping",
    top_category %in% entertainment_recreation ~ "Entertainment/Recreation",
    top_category %in% personal_services ~ "Personal Services",
    top_category %in% hospitality_lodging ~ "Hospitality/Lodging",
    top_category %in% office_professional ~ "Office/Professional",
    top_category %in% target_categories ~ "Restaurant/Bar",
    TRUE ~ "Other"
  )) |>
  filter(category_group != "Other") |> # Exclude any unintended categories
  group_by(category_group, age_group) |>
  summarize(total_visitors = sum(visitor_count, na.rm = TRUE), .groups = "drop")

# Create stacked bar plot
ggplot(df_filtered, aes(x = category_group, y = total_visitors, fill = age_group)) +
  geom_col(position = "stack") +
  labs(
    title = "Visitor Counts by POI Category and Age Group",
    x = "POI Category",
    y = "Total Visitors",
    fill = "Age Group"
  ) +
  theme_minimal() +
```

```
theme(
  axis.text.x = element_text(angle = 45, hjust = 1), # Rotate x-axis labels for readability
  legend.position = "right" # Keep legend for age group colors
)
```



Bar Plot

```
age_group_summary = df_long_model_filtered_1 |>
  group_by(age_group, top_category) |>
  summarize(total_visitors = sum(visitor_count), .groups = "drop")

ggplot(age_group_summary, aes(x = age_group, y = total_visitors, fill = top_category)) +
  geom_col(position = "dodge") +
  labs(
    title = "Visitor Counts by Age Group and Location Type",
    x = "Age Group",
    y = "Total Visitors",
    fill = "Location Type"
  ) +
  theme_minimal()
```

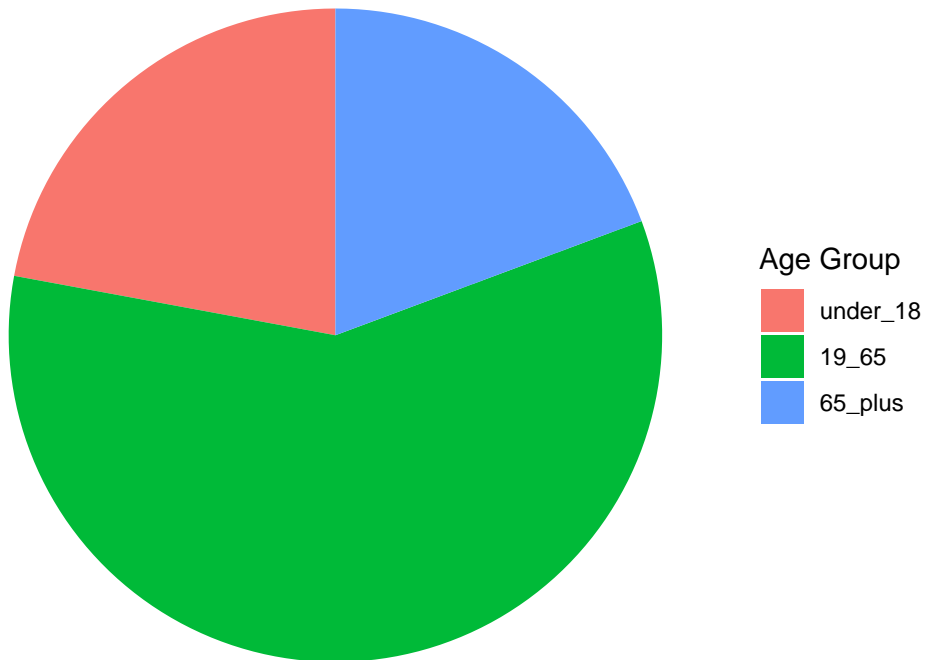
pository Credit Intermediation	Legal Services
inking Places (Alcoholic Beverages)	Machinery, Equipment, and Supplies
ycleaning and Laundry Services	Museums, Historical Sites, and Simila
ectronic and Precision Equipment Repair and Maintenance	Offices of Dentists
ectronics and Appliance Stores	Offices of Other Health Practitioners
ementary and Secondary Schools	Offices of Physicians
rists	Offices of Real Estate Agents and Br
rniture Stores	Other Amusement and Recreation In
soline Stations	Other Financial Investment Activities
neral Medical and Surgical Hospitals	Other Miscellaneous Manufacturing
neral Merchandise Stores, including Warehouse Clubs and Supercenters	Other Miscellaneous Store Retailers
ass and Glass Product Manufacturing	Other Personal Services
ocery Stores	Other Professional, Scientific, and Te
alth and Personal Care Stores	Other Schools and Instruction
me Furnishings Stores	Other Specialty Trade Contractors
restigation and Security Services	Personal and Household Goods Repa
welry, Luggage, and Leather Goods Stores	Personal Care Services
stice, Public Order, and Safety Activities	Printing and Related Support Activitie

Pie chart (alternative to above)

```
age_group_proportions = df_long_model_filtered_1 |>
  group_by(age_group) |>
  summarize(total_visitors = sum(visitor_count), .groups = "drop") |>
  mutate(proportion = total_visitors / sum(total_visitors))

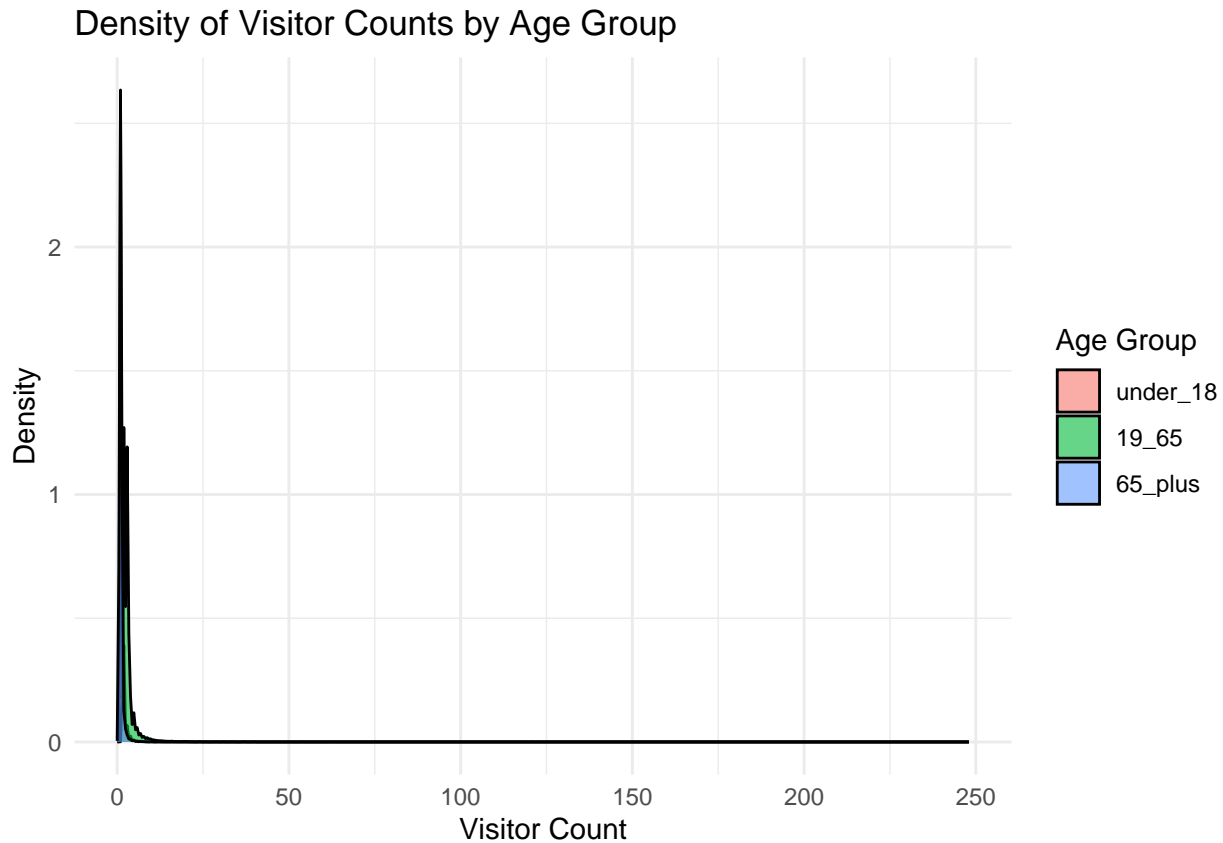
ggplot(age_group_proportions, aes(x = "", y = proportion, fill = age_group)) +
  geom_bar(stat = "identity", width = 1) +
  coord_polar(theta = "y") +
  labs(
    title = "Proportion of Visitors by Age Group",
    fill = "Age Group"
  ) +
  theme_void()
```


Proportion of Visitors by Age Group



Density Plot

```
ggplot(df_long_model_filtered_1, aes(x = visitor_count, fill = age_group)) +  
  geom_density(alpha = 0.6) +  
  labs(  
    title = "Density of Visitor Counts by Age Group",  
    x = "Visitor Count",  
    y = "Density",  
    fill = "Age Group"  
  ) +  
  theme_minimal()
```



Modeling

All categories vs. categories of interest

```
poisson_model_interact_1 = glm(visitor_count ~ age_group * non_restaurant, family = poisson(link = "log"),
summary(poisson_model_interact_1)
```

```
##
## Call:
## glm(formula = visitor_count ~ age_group * non_restaurant, family = poisson(link = "log"),
##      data = df_long_model_filtered_1)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.16988    0.02238   7.592 3.16e-14 ***
## age_group19_65    0.96487    0.02630  36.691 < 2e-16 ***
## age_group65_plus  -0.10608    0.03252  -3.262  0.00111 **
## non_restaurantYes  0.06969    0.02323   3.000  0.00270 **
## age_group19_65:non_restaurantYes  0.01149    0.02730   0.421  0.67383
## age_group65_plus:non_restaurantYes -0.03016    0.03378  -0.893  0.37195
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 85790  on 65684  degrees of freedom
```

```
## Residual deviance: 52109  on 65679  degrees of freedom
## AIC: 207746
##
## Number of Fisher Scoring iterations: 5
dispersion_test = sum(residuals(poisson_model_interact_1, type = "pearson")^2) / poisson_model_interact_1$deviance
print(dispersion_test)

## [1] 2.619864
#Overdispersion present, use NB
```

NB models, overdispersion was present

NB model on whole data

“Are older individuals visiting restaurants/bars at lower rates compared to other age groups?” A negative estimate implies that an age group is visiting a location at a lower rate than the reference

```
nb_whole = glm.nb(visitor_count ~ age_group, data = df_long)
summary(nb_whole)

##
## Call:
## glm.nb(formula = visitor_count ~ age_group, data = df_long, init.theta = 9.54321176,
##      link = log)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.234380   0.006397   36.64  <2e-16 ***
## age_group19_65  0.975531   0.007702  126.66  <2e-16 ***
## age_group65_plus -0.134037   0.009328  -14.37  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(9.5432) family taken to be 1)
##
##      Null deviance: 60126  on 65684  degrees of freedom
## Residual deviance: 32394  on 65682  degrees of freedom
## AIC: 199382
##
## Number of Fisher Scoring iterations: 1
##
##              Theta:  9.543
##             Std. Err.:  0.170
##
##      2 x log-likelihood:  -199373.606
df_model_filtered = df_long |>
  filter(top_category %in% target_categories)

nb_rest = glm.nb(visitor_count ~ age_group, data = df_model_filtered)
summary(nb_rest)

##
## Call:
```

```
## glm.nb(formula = visitor_count ~ age_group, data = df_model_filtered,
##       init.theta = 45.18341845, link = log)
##
## Coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.16988    0.02267   7.494 6.68e-14 ***
## age_group19_65  0.96487    0.02679  36.013 < 2e-16 ***
## age_group65_plus -0.10608    0.03292  -3.222  0.00127 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(45.1834) family taken to be 1)
##
## Null deviance: 4199.5  on 5054  degrees of freedom
## Residual deviance: 1982.1  on 5052  degrees of freedom
## AIC: 14056
##
## Number of Fisher Scoring iterations: 1
##
##               Theta:  45.18
##             Std. Err.:  9.56
##
## 2 x log-likelihood:  -14048.09
```

NB with interaction

“Are older individuals visiting restaurants/bars at lower rates compared to other location types?”

```
run_nb_model = function(df, category_name, category_vector, reference_name, target_categories) {
  df_model = df |>
    filter(top_category %in% c(target_categories, category_vector)) |> # Filter to only relevant POIs
    mutate(category_indicator = if_else(top_category %in% target_categories, reference_name, category_name),
           category_indicator = factor(category_indicator, levels = c(reference_name, category_name)),
           age_group = factor(age_group, levels = c("under_18", "19_65", "65_plus"))) # Ensure Restaurant/Bar

  nb_model = glm.nb(visitor_count ~ age_group * category_indicator + offset(log(total_visitors)), data = df_model)

  print(summary(nb_model))

  return(nb_model)
}

# All groups v. restaurant and bar
nb_model_1 = run_nb_model(df_long, "Non-Restaurant", all_cat, "Restaurant/Bar", target_categories)
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =  
## control$trace > : iteration limit reached  
## Warning in glm.nb(visitor_count ~ age_group * category_indicator +  
## offset(log(total_visitors))), : alternation limit reached  
  
##  
## Call:  
## glm.nb(formula = visitor_count ~ age_group * category_indicator +  
##       offset(log(total_visitors)), data = df_model, init.theta = 267928.2445,  
##       link = log)  
##  
## Coefficients:
```



```
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in glm.nb(visitor_count ~ age_group * category_indicator +
## offset(log(total_visitors)), : alternation limit reached
##
## Call:
## glm.nb(formula = visitor_count ~ age_group * category_indicator +
## offset(log(total_visitors)), data = df_model, init.theta = 416107.5965,
## link = log)
##
## Coefficients:
##
## Estimate Std. Error z value
## (Intercept) -1.509349 0.022378 -67.449
## age_group19_65 0.964865 0.026297 36.690
## age_group65_plus -0.106084 0.032520 -3.262
## category_indicatorHealthcare 0.009823 0.026047 0.377
## age_group19_65:category_indicatorHealthcare -0.020681 0.030633 -0.675
## age_group65_plus:category_indicatorHealthcare 0.010507 0.037825 0.278
## Pr(>|z|)
## (Intercept) < 2e-16 ***
## age_group19_65 < 2e-16 ***
## age_group65_plus 0.00111 **
```



```

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in glm.nb(visitor_count ~ age_group * category_indicator +
## offset(log(total_visitors)), : alternation limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group * category_indicator +
##       offset(log(total_visitors)), data = df_model, init.theta = 395592.006,
##       link = log)
##
## Coefficients:
##
##               Estimate Std. Error
## (Intercept)      -1.509349    0.022378
## age_group19_65      0.964865    0.026297
## age_group65_plus    -0.106084    0.032520
## category_indicatorEssential Services -0.006683    0.030528
## age_group19_65:category_indicatorEssential Services  0.001630    0.035872
## age_group65_plus:category_indicatorEssential Services 0.028554    0.044213
##
##               z value Pr(>|z|)
## (Intercept)      -67.449 < 2e-16 ***
## age_group19_65     36.690 < 2e-16 ***
## age_group65_plus    -3.262  0.00111 **
## category_indicatorEssential Services   -0.219  0.82673
## age_group19_65:category_indicatorEssential Services   0.045  0.96376
## age_group65_plus:category_indicatorEssential Services   0.646  0.51839
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(395581.6) family taken to be 1)
##
## Null deviance: 5922.09  on 11024  degrees of freedom
## Residual deviance:  976.63  on 11019  degrees of freedom
## AIC: 26794

```

```

##
## Number of Fisher Scoring iterations: 1
##
##
##          Theta: 395592
##          Std. Err.: 1111370
## Warning while fitting theta: alternation limit reached
##
## 2 x log-likelihood: -26780.16
# Retail shopping v. restaurant and bar
nb_model_4 = run_nb_model(df_long, "Retail Shopping", retail_shopping, "Restaurant/Bar", target_category)

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group * category_indicator +
##       offset(log(total_visitors)), data = df_model, init.theta = 320795.0344,
##       link = log)
##
## Coefficients:
##
##              Estimate Std. Error z value
## (Intercept)      -1.50935    0.02238 -67.449
## age_group19_65       0.96486    0.02630  36.690
## age_group65_plus    -0.10608    0.03252  -3.262
## category_indicatorRetail Shopping    0.02902    0.02585   1.123
## age_group19_65:category_indicatorRetail Shopping -0.01954    0.03040  -0.643
## age_group65_plus:category_indicatorRetail Shopping -0.09148    0.03780  -2.420
##
##              Pr(>|z|)
## (Intercept)      < 2e-16 ***
## age_group19_65    < 2e-16 ***
## age_group65_plus   0.00111 **
## category_indicatorRetail Shopping    0.26150
## age_group19_65:category_indicatorRetail Shopping    0.52031
## age_group65_plus:category_indicatorRetail Shopping    0.01551 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(320795) family taken to be 1)
##
##      Null deviance: 11434.5  on 18149  degrees of freedom
## Residual deviance:  2033.2  on 18144  degrees of freedom
## AIC: 45208
##
## Number of Fisher Scoring iterations: 1
##
##
##          Theta: 320795
##          Std. Err.: 631190
## Warning while fitting theta: iteration limit reached
##

```

```
## 2 x log-likelihood: -45194.49
# Entertainment/Recreation v. restaurant and bar
nb_model_5 = run_nb_model(df_long, "Entertainment/Recreation", entertainment_recreation, "Restaurant/Bar")

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group * category_indicator +
## offset(log(total_visitors)), data = df_model, init.theta = 411283.3029,
## link = log)
##
## Coefficients:
##
## Estimate Std. Error
## (Intercept) -1.50935 0.02238
## age_group19_65 0.96486 0.02630
## age_group65_plus -0.10608 0.03252
## category_indicatorEntertainment/Recreation 0.03222 0.04157
## age_group19_65:category_indicatorEntertainment/Recreation -0.05029 0.04909
## age_group65_plus:category_indicatorEntertainment/Recreation -0.01646 0.06059
##
## z value Pr(>|z|)
## (Intercept) -67.449 < 2e-16
## age_group19_65 36.690 < 2e-16
## age_group65_plus -3.262 0.00111
## category_indicatorEntertainment/Recreation 0.775 0.43831
## age_group19_65:category_indicatorEntertainment/Recreation -1.024 0.30562
## age_group65_plus:category_indicatorEntertainment/Recreation -0.272 0.78584
##
## (Intercept) ***
## age_group19_65 ***
## age_group65_plus **
## category_indicatorEntertainment/Recreation
## age_group19_65:category_indicatorEntertainment/Recreation
## age_group65_plus:category_indicatorEntertainment/Recreation
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(411283.3) family taken to be 1)
##
## Null deviance: 3749.83 on 7130 degrees of freedom
## Residual deviance: 597.25 on 7125 degrees of freedom
## AIC: 17310
##
## Number of Fisher Scoring iterations: 1
##
##
## Theta: 411283
## Std. Err.: 1433299
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -17295.6
```

```
# Personal Services v. restaurant and bar
```

```
nb_model_6 = run_nb_model(df_long, "Personal Services", personal_services, "Restaurant/Bar", target_cat
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
```

```
## control$trace > : iteration limit reached
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
```

```
## control$trace > : iteration limit reached
```

```
##
```

```
## Call:
```

```
## glm.nb(formula = visitor_count ~ age_group * category_indicator +
```

```
##   offset(log(total_visitors)), data = df_model, init.theta = 209777.2543,
```

```
##   link = log)
```

```
##
```

```
## Coefficients:
```

```
##                                     Estimate Std. Error
```

```
## (Intercept)                       -1.50935      0.02238
```

```
## age_group19_65                     0.96486      0.02630
```

```
## age_group65_plus                   -0.10608      0.03252
```

```
## category_indicatorPersonal Services  0.00100      0.03308
```

```
## age_group19_65:category_indicatorPersonal Services  0.03614      0.03877
```

```
## age_group65_plus:category_indicatorPersonal Services -0.11929      0.04893
```

```
##                                     z value Pr(>|z|)
```

```
## (Intercept)                       -67.449 < 2e-16 ***
```

```
## age_group19_65                     36.690 < 2e-16 ***
```

```
## age_group65_plus                   -3.262 0.00111 **
```

```
## category_indicatorPersonal Services  0.030 0.97588
```

```
## age_group19_65:category_indicatorPersonal Services  0.932 0.35125
```

```
## age_group65_plus:category_indicatorPersonal Services -2.438 0.01478 *
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## (Dispersion parameter for Negative Binomial(209777.3) family taken to be 1)
```

```
##
```

```
##   Null deviance: 5634.25  on 8348  degrees of freedom
```

```
## Residual deviance:  990.19  on 8343  degrees of freedom
```

```
## AIC: 21037
```

```
##
```

```
## Number of Fisher Scoring iterations: 1
```

```
##
```

```
##
```

```
##               Theta: 209777
```

```
##           Std. Err.: 794709
```

```
## Warning while fitting theta: iteration limit reached
```

```
##
```

```
## 2 x log-likelihood: -21023.24
```

```
#Hospitality/Lodging v. restaurant and bar
```

```
nb_model_7 = run_nb_model(df_long, "Hospitality/Lodging", hospitality_lodging, "Restaurant/Bar", target_cat
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
```

```
## control$trace > : iteration limit reached
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
```

```
## control$trace > : iteration limit reached
```

```
##
## Call:
## glm.nb(formula = visitor_count ~ age_group * category_indicator +
##       offset(log(total_visitors)), data = df_model, init.theta = 189020.2782,
##       link = log)
##
## Coefficients:
##
##               Estimate Std. Error
## (Intercept)      -1.509348   0.022378
## age_group19_65      0.964863   0.026298
## age_group65_plus    -0.106084   0.032520
## category_indicatorHospitality/Lodging  0.004875   0.041567
## age_group19_65:category_indicatorHospitality/Lodging  0.076046   0.048496
## age_group65_plus:category_indicatorHospitality/Lodging -0.294487   0.064157
##
##               z value Pr(>|z|)
## (Intercept)      -67.449 < 2e-16 ***
## age_group19_65     36.690 < 2e-16 ***
## age_group65_plus    -3.262  0.00111 **
## category_indicatorHospitality/Lodging  0.117  0.90663
## age_group19_65:category_indicatorHospitality/Lodging  1.568  0.11686
## age_group65_plus:category_indicatorHospitality/Lodging -4.590 4.43e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(189020.3) family taken to be 1)
##
##      Null deviance: 4488.25  on 6320  degrees of freedom
## Residual deviance:  787.67  on 6315  degrees of freedom
## AIC: 16022
##
## Number of Fisher Scoring iterations: 1
##
##
##           Theta: 189020
##          Std. Err.: 865917
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -16008.01
```

```
#Office/Professional v. restaurant and bar
```

```
nb_model_8 = run_nb_model(df_long, "Office/Professional", office_professional, "Restaurant/Bar", target,
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
```

```
##
## Call:
## glm.nb(formula = visitor_count ~ age_group * category_indicator +
##       offset(log(total_visitors)), data = df_model, init.theta = 372297.1849,
##       link = log)
##
## Coefficients:
##
##               Estimate Std. Error
```

```

## (Intercept) -1.509349 0.022378
## age_group19_65 0.964865 0.026298
## age_group65_plus -0.106084 0.032520
## category_indicatorOffice/Professional 0.003127 0.025891
## age_group19_65:category_indicatorOffice/Professional -0.008028 0.030435
## age_group65_plus:category_indicatorOffice/Professional 0.007597 0.037607
## z value Pr(>|z|)
## (Intercept) -67.449 < 2e-16 ***
## age_group19_65 36.690 < 2e-16 ***
## age_group65_plus -3.262 0.00111 **
## category_indicatorOffice/Professional 0.121 0.90387
## age_group19_65:category_indicatorOffice/Professional -0.264 0.79196
## age_group65_plus:category_indicatorOffice/Professional 0.202 0.83992
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(372297.2) family taken to be 1)
##
## Null deviance: 10670.2 on 19712 degrees of freedom
## Residual deviance: 1712.3 on 19707 degrees of freedom
## AIC: 47947
##
## Number of Fisher Scoring iterations: 1
##
##
## Theta: 372297
## Std. Err.: 827312
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -47932.69

extract_nb_results = function(model, category_name) {
  results = broom.mixed::tidy(model) |>
    filter(grepl("category_indicator", term)) |>
    mutate(category = category_name) |>
    relocate(category) |>
    mutate(significance = case_when(
      p.value < 0.001 ~ "***",
      p.value < 0.01 ~ "**",
      p.value < 0.05 ~ "*",
      TRUE ~ ""
    ))

  return(results)
}

nb_summary_table = bind_rows(
  extract_nb_results(nb_model_1, "Non-Restaurant"),
  extract_nb_results(nb_model_2, "Healthcare"),
  extract_nb_results(nb_model_3, "Essential Services"),
  extract_nb_results(nb_model_4, "Retail Shopping"),
  extract_nb_results(nb_model_5, "Entertainment/Recreation"),
  extract_nb_results(nb_model_6, "Personal Services"),
  extract_nb_results(nb_model_7, "Hospitality/Lodging"),

```

```

extract_nb_results(nb_model_8, "Office/Professional")
)

knitr::kable(nb_summary_table)

```

category	term	estimate	std.error	statistic	p.value	significance
Non-Restaurant	category_indicatorNon-Restaurant	- 0.0232313	- 0.9726820			
		0.0007955	0.0342447			
Non-Restaurant	age_group19_65:category_indicatorNon-Restaurant	0.0114845	0.0272978	0.4207103	0.6739667	
Non-Restaurant	age_group65_plus:category_indicatorNon-Restaurant	- 0.0337804	- 0.3719488			
		0.0301601	0.8928290			
Healthcare	category_indicatorHealthcare	0.0098230	0.0260468	0.3771287	0.7060779	
Healthcare	age_group19_65:category_indicatorHealthcare	- 0.0306327	- 0.4995838			
		0.0206815	0.6751448			
Healthcare	age_group65_plus:category_indicatorHealthcare	0.0105075	0.0378251	0.2777916	0.7811724	
Essential Services	category_indicatorEssential Services	- 0.0305283	- 0.8267289			
		0.0066826	0.2188988			
Essential Services	age_group19_65:category_indicatorEssential Services	0.0016300	0.0358723	0.0454397	0.9637569	
Essential Services	age_group65_plus:category_indicatorEssential Services	0.0285537	0.0442126	0.6458282	0.5183906	
Retail Shopping	category_indicatorRetail Shopping	0.0290238	0.0258485	1.1228414	0.2615049	
Retail Shopping	age_group19_65:category_indicatorRetail Shopping	- 0.0303973	- 0.5203116			
		0.0195414	0.6428651			
Retail Shopping	age_group65_plus:category_indicatorRetail Shopping	- 0.0378008	- 0.0155132*			
		0.0914844	2.4201720			
Entertainment/Recreation	category_indicatorEntertainment/Recreation	0.0322161	0.0415662	0.7750540	0.4383078	
Entertainment/Recreation	age_group19_65:category_indicatorEntertainment/Recreation	- 0.0499938	- 0.3056216			
		0.0502943	1.0244525			
Entertainment/Recreation	age_group65_plus:category_indicatorEntertainment/Recreation	0.0605930	- 0.7858370			
		0.0164644	0.2717204			
Personal Services	category_indicatorPersonal Services	0.0010003	0.0330793	0.0302387	0.9758767	
Personal Services	age_group19_65:category_indicatorPersonal Services	0.0361410	0.0387706	0.9321740	0.3512466	
Personal Services	age_group65_plus:category_indicatorPersonal Services	- 0.0489341	- 0.0147809*			
		0.1192867	2.4377026			
Hospitality/Lodging	category_indicatorHospitality/Lodging	0.0048754	0.0415667	0.1172917	0.9066289	
Hospitality/Lodging	age_group19_65:category_indicatorHospitality/Lodging	0.0760165	0.0484965	1.5680823	0.1168619	
Hospitality/Lodging	age_group65_plus:category_indicatorHospitality/Lodging	- 0.0641570	- 0.0000044***			
		0.2944870	4.5900955			
Office/Professional	category_indicatorOffice/Professional	0.0031269	0.0258913	0.1207685	0.9038744	
Office/Professional	age_group19_65:category_indicatorOffice/Professional	0.0304354	- 0.7919613			
		0.0080278	0.2637647			
Office/Professional	age_group65_plus:category_indicatorOffice/Professional	0.0750650	0.0376074	0.2020005	0.8399163	

NB no interaction and offset

```

run_nb_model = function(df, category_name, category_vector) {
  df_model = df |>
    filter(top_category %in% category_vector) |> # Filter to only relevant POIs

```

```

    mutate(
      age_group = factor(age_group, levels = c("under_18", "19_65", "65_plus"))) # Ensure age gro

nb_model = glm.nb(visitor_count ~ age_group + offset(log(total_visitors)), data = df_model)

print(summary(nb_model))

return(nb_model)
}

# Run models for each POI type vs. Restaurant/Bar
nb_model_1 = run_nb_model(df_long, "Full", all_cat)

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
## data = df_model, init.theta = 267771.141, link = log)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -1.510087   0.006011 -251.23  <2e-16 ***
## age_group19_65   0.975524   0.007053  138.31  <2e-16 ***
## age_group65_plus -0.134037   0.008800  -15.23  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(267771.1) family taken to be 1)
##
## Null deviance: 41075.8 on 65684 degrees of freedom
## Residual deviance: 7440.3 on 65682 degrees of freedom
## AIC: 163074
##
## Number of Fisher Scoring iterations: 1
##
##
##              Theta: 267771
##              Std. Err.: 327730
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -163065.8
nb_model_2 = run_nb_model(df_long, "Healthcare", medical_services)

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =

```



```
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in glm.nb(visitor_count ~ age_group + offset(log(total_visitors)), :
## alternation limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
##       data = df_model, init.theta = 421352.8943, link = log)
##
## Coefficients:
```



```

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in glm.nb(visitor_count ~ age_group + offset(log(total_visitors)), :
## alternation limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
## data = df_model, init.theta = 389717.6916, link = log)
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.51603    0.02077  -73.006 < 2e-16 ***
## age_group19_65    0.96649    0.02440   39.614 < 2e-16 ***
## age_group65_plus -0.07753    0.02995   -2.588  0.00964 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(389747.2) family taken to be 1)
##
## Null deviance: 3179.16 on 5969 degrees of freedom
## Residual deviance: 545.16 on 5967 degrees of freedom
## AIC: 14484
##
## Number of Fisher Scoring iterations: 1
##
##
##             Theta: 389718
##             Std. Err.: 1506501
## Warning while fitting theta: alternation limit reached
##
## 2 x log-likelihood: -14476.38

```

```

nb_model_4 = run_nb_model(df_long, "Retail Shopping", retail_shopping)

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
## data = df_model, init.theta = 297023.5468, link = log)
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.48032    0.01294 -114.42  <2e-16 ***
## age_group19_65  0.94532    0.01525  62.01  <2e-16 ***
## age_group65_plus -0.19757    0.01927 -10.25  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(297023.5) family taken to be 1)
##
## Null deviance: 8691.5 on 13094 degrees of freedom
## Residual deviance: 1601.8 on 13092 degrees of freedom
## AIC: 32899
##
## Number of Fisher Scoring iterations: 1
##
##
##             Theta: 297024
##           Std. Err.: 667423
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -32890.71

nb_model_5 = run_nb_model(df_long, "Entertainment/Recreation", entertainment_recreation)

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
## data = df_model, init.theta = 434611.9701, link = log)
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.47713    0.03503 -42.169  <2e-16 ***
## age_group19_65  0.91457    0.04146  22.061  <2e-16 ***
## age_group65_plus -0.12255    0.05113  -2.397   0.0165 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##

```

```

## (Dispersion parameter for Negative Binomial(434612) family taken to be 1)
##
## Null deviance: 1006.91 on 2075 degrees of freedom
## Residual deviance: 165.79 on 2073 degrees of freedom
## AIC: 4999.8
##
## Number of Fisher Scoring iterations: 1
##
##
## Theta: 434612
## Std. Err.: 2890549
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -4991.82
nb_model_6 = run_nb_model(df_long, "Personal Services", personal_services)

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
## data = df_model, init.theta = 499.4942562, link = log)
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.50589 0.02443 -61.640 < 2e-16 ***
## age_group19_65 0.99395 0.02864 34.700 < 2e-16 ***
## age_group65_plus -0.22499 0.03665 -6.139 8.3e-10 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(499.4943) family taken to be 1)
##
## Null deviance: 2814.63 on 3293 degrees of freedom
## Residual deviance: 541.26 on 3291 degrees of freedom
## AIC: 8725
##
## Number of Fisher Scoring iterations: 1
##
##
## Theta: 499
## Std. Err.: 357
##
## 2 x log-likelihood: -8717.025
nb_model_7 = run_nb_model(df_long, "Hospitality/Lodging", hospitality_lodging)

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
## data = df_model, init.theta = 134.8428159, link = log)
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.49051 0.03543 -42.066 < 2e-16 ***
## age_group19_65 1.00909 0.04165 24.230 < 2e-16 ***

```



```

## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
## Warning in glm.nb(visitor_count ~ age_group + offset(log(total_visitors)), :
## alternation limit reached

##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
##       data = df_model, init.theta = 360608.9225, link = log)
##
## Coefficients:
##              Estimate Std. Error  z value Pr(>|z|)
## (Intercept)    -1.50622    0.01302  -115.656 < 2e-16 ***
## age_group19_65    0.95684    0.01532   62.450 < 2e-16 ***
## age_group65_plus -0.09849    0.01889   -5.214 1.85e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(360608.7) family taken to be 1)
##
##      Null deviance: 7927.3  on 14657  degrees of freedom
## Residual deviance: 1280.8  on 14655  degrees of freedom
## AIC: 35637
##
## Number of Fisher Scoring iterations: 1
##
##
##              Theta: 360609
##            Std. Err.: 955628
## Warning while fitting theta: alternation limit reached
##
## 2 x log-likelihood: -35628.91
nb_model_9 = run_nb_model(df_long, "Restaurant/Bar", target_categories)

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached

```

```
##
## Call:
## glm.nb(formula = visitor_count ~ age_group + offset(log(total_visitors)),
##       data = df_model, init.theta = 402384.7886, link = log)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.50935    0.02238 -67.449 < 2e-16 ***
## age_group19_65  0.96486    0.02630  36.690 < 2e-16 ***
## age_group65_plus -0.10608    0.03252  -3.262  0.00111 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(402384.8) family taken to be 1)
##
## Null deviance: 2742.92 on 5054 degrees of freedom
## Residual deviance: 431.46 on 5052 degrees of freedom
## AIC: 12312
##
## Number of Fisher Scoring iterations: 1
##
##              Theta: 402385
##              Std. Err.: 1646220
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -12303.78

extract_nb_results = function(model, category_name) {
  results = broom.mixed::tidy(model) |>
    mutate(significance = case_when(
      p.value < 0.001 ~ "***",
      p.value < 0.01 ~ "**",
      p.value < 0.05 ~ "*",
      TRUE ~ ""
    ))

  return(results)
}

# Combine results from all models
nb_summary_table = bind_rows(
  extract_nb_results(nb_model_1, "Non-Restaurant"),
  extract_nb_results(nb_model_2, "Healthcare"),
  extract_nb_results(nb_model_3, "Essential Services"),
  extract_nb_results(nb_model_4, "Retail Shopping"),
  extract_nb_results(nb_model_5, "Entertainment/Recreation"),
  extract_nb_results(nb_model_6, "Personal Services"),
  extract_nb_results(nb_model_7, "Hospitality/Lodging"),
  extract_nb_results(nb_model_8, "Office/Professional"),
  extract_nb_results(nb_model_9, "Restaurant/Bar")
)

nb_summary_table = nb_summary_table |>
```



```

bind_cols(category = c("Full","Full","Full", "Healthcare", "Healthcare", "Healthcare", "Essential Ser
relocate(category)

# Display as a table
knitr::kable(nb_summary_table)

```

category	term	estimate	std.error	statistic	p.value	significance
Full	(Intercept)	-	0.0060108	-	0.0000000	***
		1.5100867		251.227448		
Full	age_group19_65	0.9755244	0.0070534	138.304611	0.0000000	***
Full	age_group65_plus	-	0.0088001	-15.231263	0.0000000	***
		0.1340373				
Healthcare	(Intercept)	-	0.0133298	-	0.0000000	***
		1.4995257		112.494161		
Healthcare	age_group19_65	0.9441834	0.0157100	60.100941	0.0000000	***
Healthcare	age_group65_plus	-	0.0193181	-4.947521	0.0000008	***
		0.0955766				
Essential Services	(Intercept)	-	0.0207659	-73.005855	0.0000000	***
		1.5160313				
Essential Services	age_group19_65	0.9664948	0.0243981	39.613564	0.0000000	***
Essential Services	age_group65_plus	-	0.0299534	-2.588363	0.0096433	**
		0.0775303				
Retail Shopping	(Intercept)	-	0.0129381	-	0.0000000	***
		1.4803247		114.416216		
Retail Shopping	age_group19_65	0.9453229	0.0152458	62.005279	0.0000000	***
Retail Shopping	age_group65_plus	-	0.0192705	-10.252381	0.0000000	***
		0.1975684				
Entertainment/Recreation	(Intercept)	-	0.0350285	-42.169406	0.0000000	***
		1.4771327				
Entertainment/Recreation	age_group19_65	0.9145706	0.0414565	22.060946	0.0000000	***
Entertainment/Recreation	age_group65_plus	-	0.0511270	-2.396944	0.0165324	*
		0.1225485				
Personal Services	(Intercept)	-	0.0244305	-61.639888	0.0000000	***
		1.5058942				
Personal Services	age_group19_65	0.9939530	0.0286446	34.699514	0.0000000	***
Personal Services	age_group65_plus	-	0.0366484	-6.139045	0.0000000	***
		0.2249864				
Hospitality/Lodging	(Intercept)	-	0.0354329	-42.065693	0.0000000	***
		1.4905110				
Hospitality/Lodging	age_group19_65	1.0090921	0.0416464	24.230017	0.0000000	***
Hospitality/Lodging	age_group65_plus	-	0.0557730	-7.193592	0.0000000	***
		0.4012080				
Office/Professional	(Intercept)	-	0.0130233	-	0.0000000	***
		1.5062218		115.655558		
Office/Professional	age_group19_65	0.9568368	0.0153217	62.449665	0.0000000	***
Office/Professional	age_group65_plus	-	0.0188883	-5.214190	0.0000002	***
		0.0984874				
Restaurant/Bar	(Intercept)	-	0.0223775	-67.449346	0.0000000	***
		1.5093487				
Restaurant/Bar	age_group19_65	0.9648648	0.0262975	36.690365	0.0000000	***
Restaurant/Bar	age_group65_plus	-	0.0325200	-3.262121	0.0011058	**
		0.1060841				