

EDA

The objective of the Exploratory Data Analysis (EDA) is to investigate the data so as to discover patterns, to spot anomalies, to handle missing values, to test hypothesis and to check assumptions with the help of summary statistics and graphical representations.

Imports

```
library(rstudioapi)
library(tidyverse)
library(magrittr)
library(DescTools)
library(gridExtra)
```

Loading data

```
current_path = rstudioapi::getActiveDocumentContext()$path
setwd(dirname(current_path))

X_train = read_csv("../data/training_set_features.csv")
Y_train = read_csv("../data/training_set_labels.csv")

df_train = X_train %>%
  inner_join(Y_train, by='respondent_id')

dimensions = dim(df_train)
sprintf("The dataset contains %d rows and %d columns", dimensions[1], dimensions[2]) %>% cat()
```

The dataset contains 26707 rows and 38 columns

```
features = names(X_train)[-1]
targets = names(Y_train)[-1]
```

```
df_train %>% glimpse()
```

```
## Rows: 26,707
## Columns: 38
## $ respondent_id      <dbl> 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 1...
## $ h1n1_concern        <dbl> 1, 3, 1, 1, 2, 3, 0, 1, 0, 2, 2, 1, 1, ...
## $ h1n1_knowledge      <dbl> 0, 2, 1, 1, 1, 1, 0, 0, 2, 1, 1, 2, 1, ...
## $ behavioral_antiviral_meds <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ behavioral_avoidance <dbl> 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, ...
## $ behavioral_face_mask <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ behavioral_wash_hands <dbl> 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, ...
## $ behavioral_large_gatherings <dbl> 0, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 1, ...
## $ behavioral_outside_home <dbl> 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, ...
## $ behavioral_touch_face <dbl> 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 1, ...
```

```
## $ doctor_recc_h1n1      <dbl> 0, 0, NA, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, ...
## $ doctor_recc_seasonal  <dbl> 0, 0, NA, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, ...
## $ chronic_med_condition <dbl> 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, ...
## $ child_under_6_months  <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, ...
## $ health_worker         <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ health_insurance      <dbl> 1, 1, NA, NA, NA, NA, NA, 1, NA, 1, 0, ...
## $ opinion_h1n1_vacc_effective <dbl> 3, 5, 3, 3, 3, 5, 4, 5, 4, 4, 4, 3, 3, ...
## $ opinion_h1n1_risk      <dbl> 1, 4, 1, 3, 3, 2, 1, 2, 1, 2, 1, 2, 2, ...
## $ opinion_h1n1_sick_from_vacc <dbl> 2, 4, 1, 5, 2, 1, 1, 1, 1, 2, 2, 2, 1, ...
## $ opinion_seas_vacc_effective <dbl> 2, 4, 4, 5, 3, 5, 4, 4, 4, 4, 5, 4, 5, ...
## $ opinion_seas_risk      <dbl> 1, 2, 1, 4, 1, 4, 2, 2, 2, 2, 4, 2, 4, ...
## $ opinion_seas_sick_from_vacc <dbl> 2, 4, 2, 1, 4, 4, 1, 1, 1, 2, 4, 1, 1, ...
## $ age_group             <chr> "55 - 64 Years", "35 - 44 Years", "18 -...
## $ education             <chr> "< 12 Years", "12 Years", "College Grad...
## $ race                  <chr> "White", "White", "White", "White", "Wh...
## $ sex                   <chr> "Female", "Male", "Male", "Female", "Fe...
## $ income_poverty        <chr> "Below Poverty", "Below Poverty", "<= $...
## $ marital_status        <chr> "Not Married", "Not Married", "Not Marr...
## $ rent_or_own           <chr> "Own", "Rent", "Own", "Rent", "Own", "O...
## $ employment_status     <chr> "Not in Labor Force", "Employed", "Empl...
## $ hhs_geo_region        <chr> "oxchjgsf", "bhuqouqj", "qufhixun", "lr...
## $ census_msa            <chr> "Non-MSA", "MSA, Not Principle City", ...
## $ household_adults      <dbl> 0, 0, 2, 0, 1, 2, 0, 2, 1, 0, 2, 1, 1, ...
## $ household_children    <dbl> 0, 0, 0, 0, 0, 3, 0, 0, 0, 0, 0, 2, 0, ...
## $ employment_industry   <chr> NA, "pxcmvdjn", "rucpziiij", NA, "wxleye...
## $ employment_occupation <chr> NA, "xgwztkwe", "xtkaffoo", NA, "emcorr...
## $ h1n1_vaccine          <dbl> 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 0, ...
## $ seasonal_vaccine      <dbl> 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, ...
```

- All of the columns available are categorical
- Some of them are labeled as numeric values and others are still characters
- The features `hhs_geo_region`, `employment_industry`, `employment_occupation` are random character strings

Targets

Checking class balance.

```
for(target in targets){
  percent = df_train[[target]] %>% mean()
  sprintf("- %.1f%% of the observations of the column %s are 1's \n", percent * 100, target) %>% cat()
}
```

- 21.2% of the observations of the column `h1n1_vaccine` are 1's
- 46.6% of the observations of the column `seasonal_vaccine` are 1's

The classes of the variable corresponding to Whether the respondent received H1N1 flu vaccine is imbalanced, much more people didn't received the H1N1 vaccine.

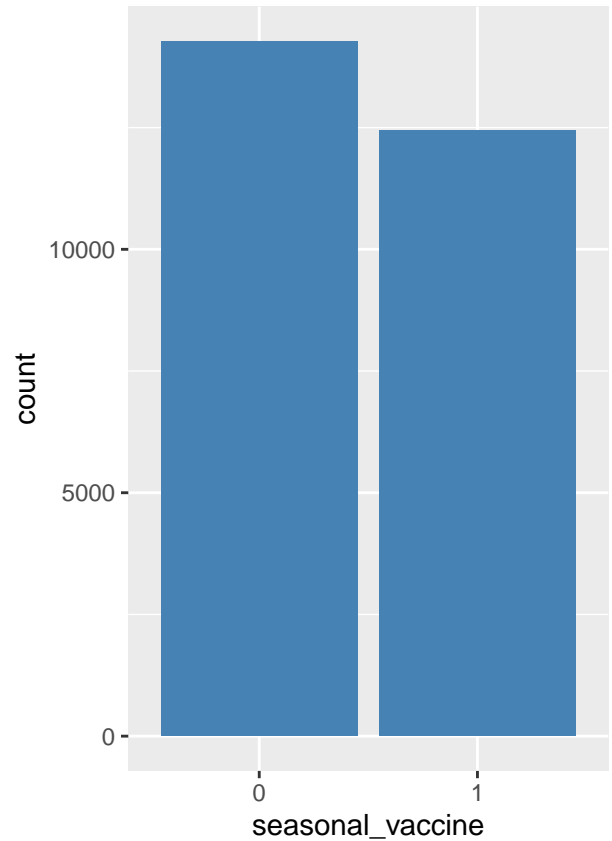
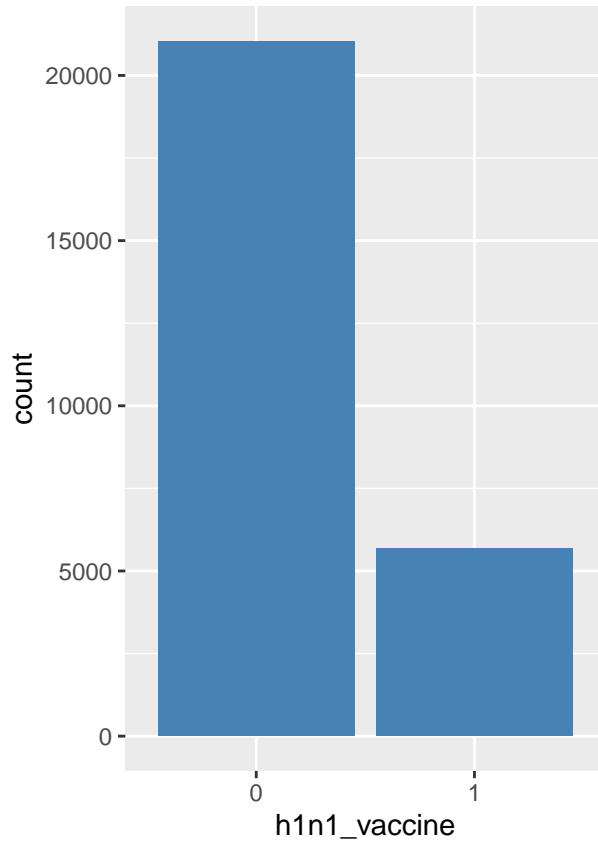
```
target = "h1n1_vaccine"

df_train = df_train %>%
  mutate_at(targets, factor)

plots = list()
i = 1
```

```
for(target in targets){
  p = df_train %>%
    ggplot(aes_string(x=target)) +
    geom_bar(fill="steelblue")
  plots[[i]] = p
  i = i + 1
}

do.call("grid.arrange", c(plots, ncol=2))
```



Features

```
df_train %>%
  select_if(is.character)
```

```
## # A tibble: 26,707 x 12
##   age_group education race sex income_poverty marital_status rent_or_own
##   <chr>      <chr>    <chr> <chr> <chr>          <chr>          <chr>
## 1 55 - 64 ~ < 12 Yea~ White Fema~ Below Poverty Not Married Own
## 2 35 - 44 ~ 12 Years White Male Below Poverty Not Married Rent
## 3 18 - 34 ~ College ~ White Male <= $75,000, A~ Not Married Own
## 4 65+ Years 12 Years White Fema~ Below Poverty Not Married Rent
## 5 45 - 54 ~ Some Col~ White Fema~ <= $75,000, A~ Married Own
## 6 65+ Years 12 Years White Male <= $75,000, A~ Married Own
## 7 55 - 64 ~ < 12 Yea~ White Male <= $75,000, A~ Not Married Own
```

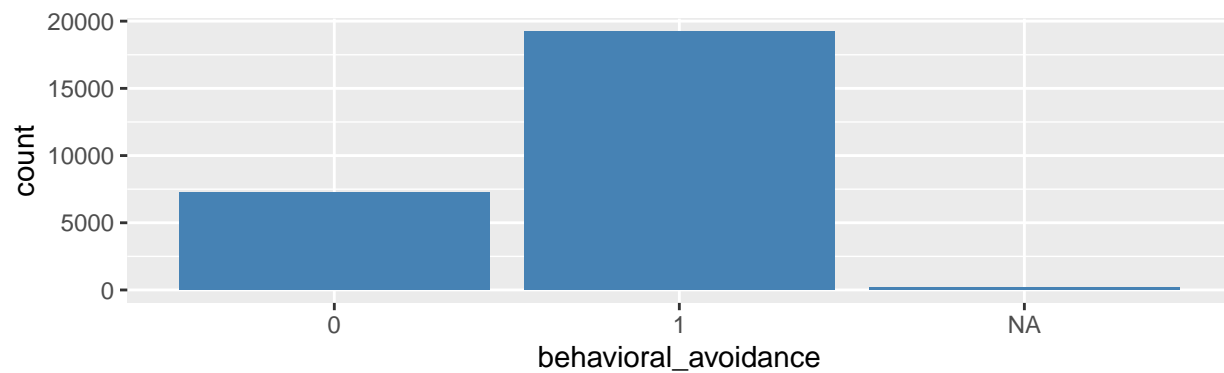
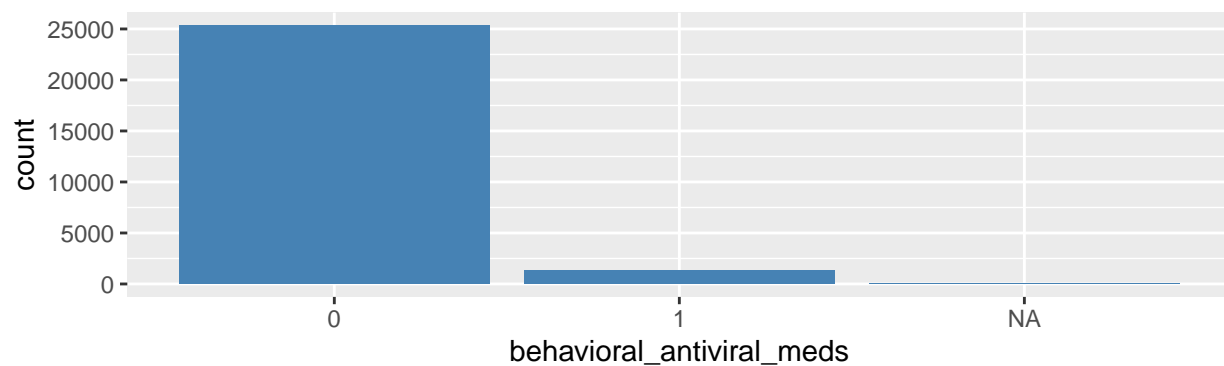
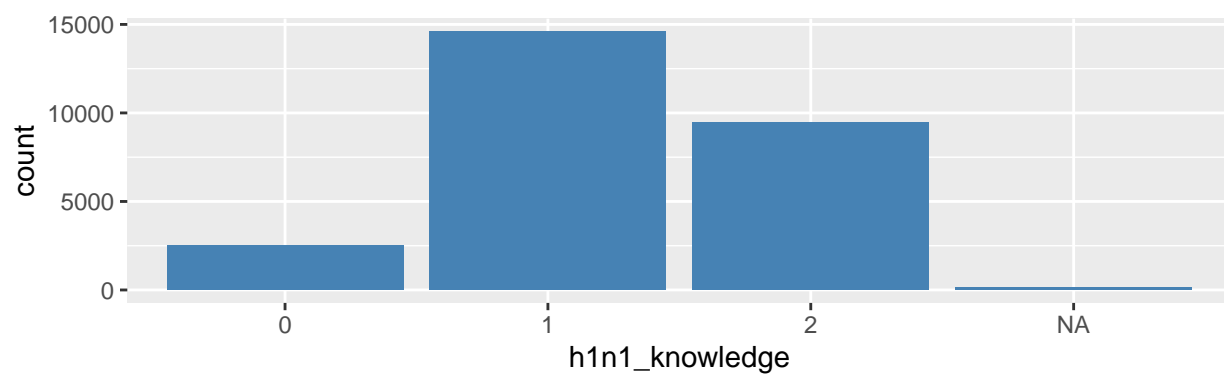
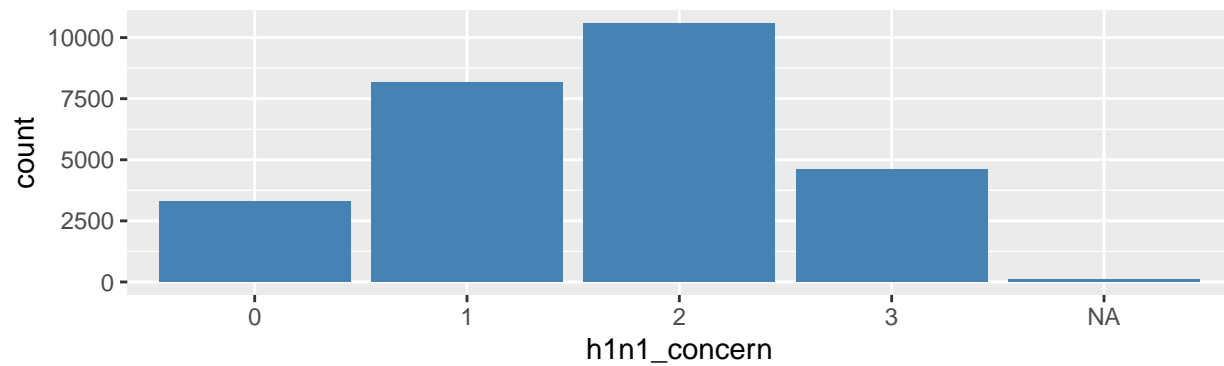
```
## 8 45 - 54 ~ Some Col~ White Fema~ <= $75,000, A~ Married      Own
## 9 45 - 54 ~ College ~ White Male > $75,000      Married      Own
## 10 55 - 64 ~ 12 Years White Male <= $75,000, A~ Not Married  Own
## # ... with 26,697 more rows, and 5 more variables: employment_status <chr>,
## #   hhs_geo_region <chr>, census_msa <chr>, employment_industry <chr>,
## #   employment_occupation <chr>
```

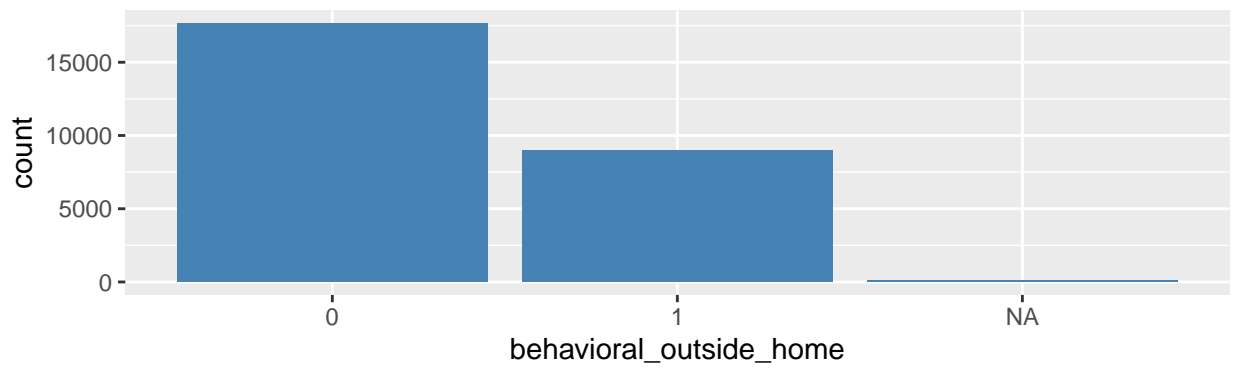
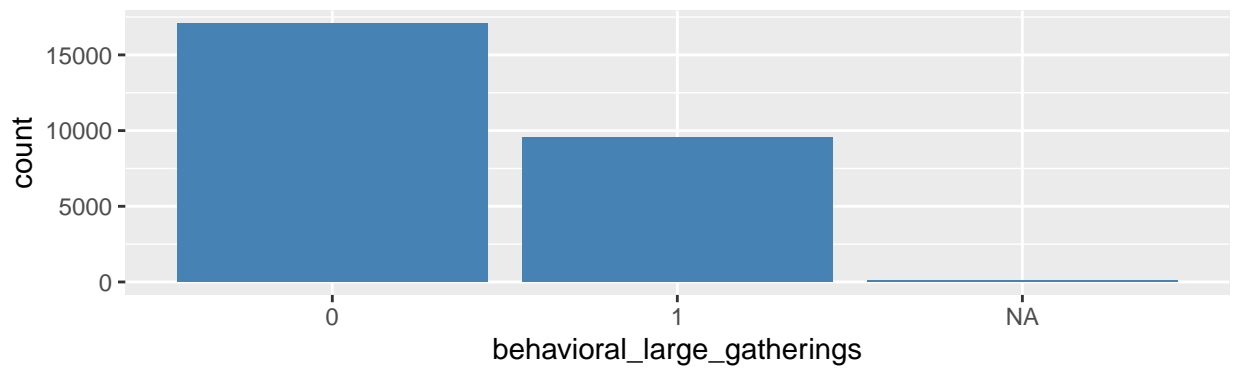
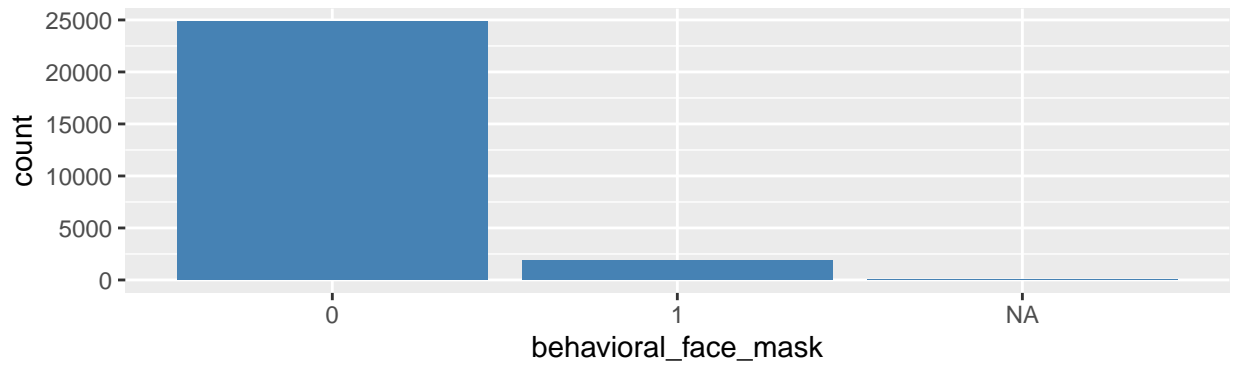
Distribution of each feature:

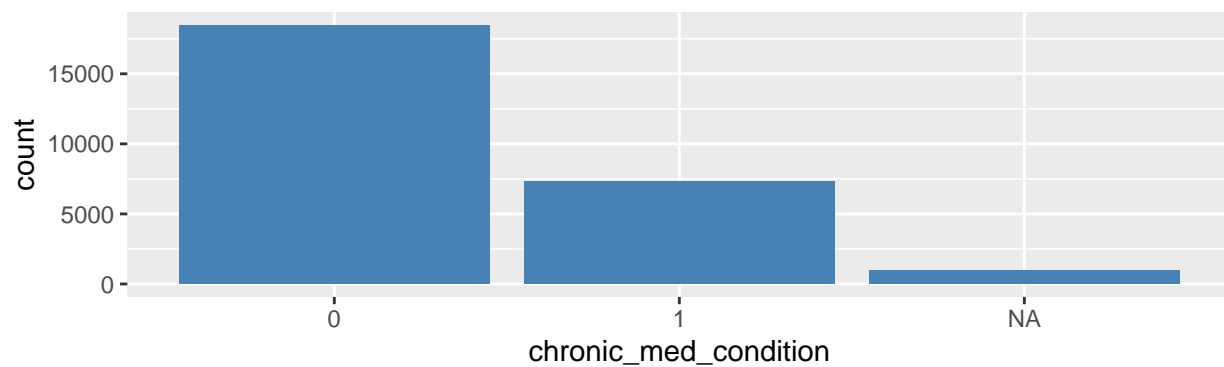
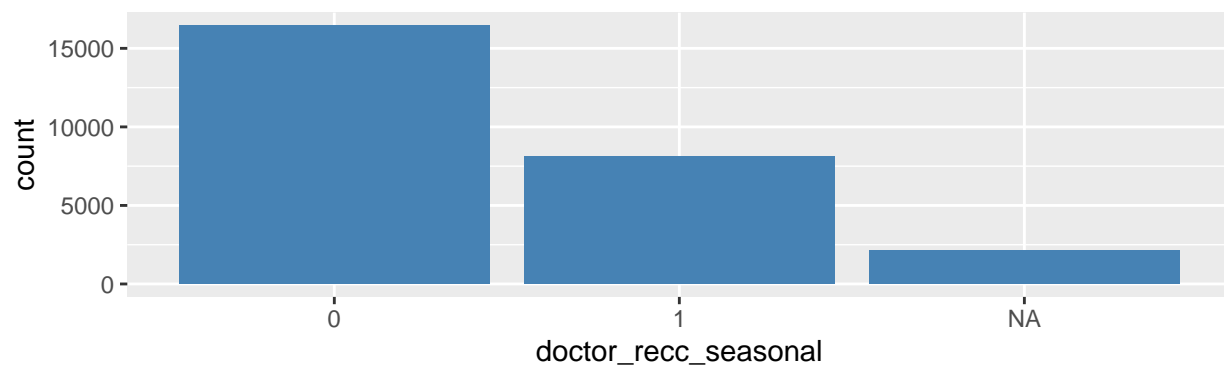
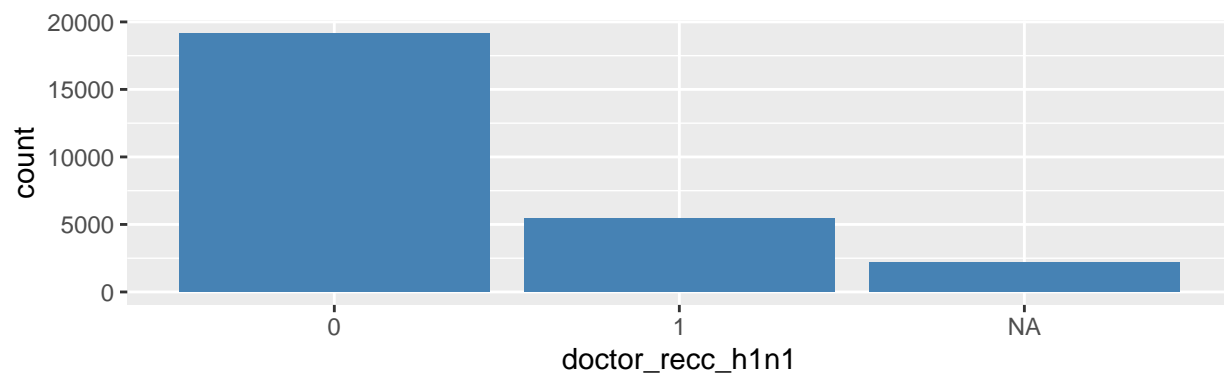
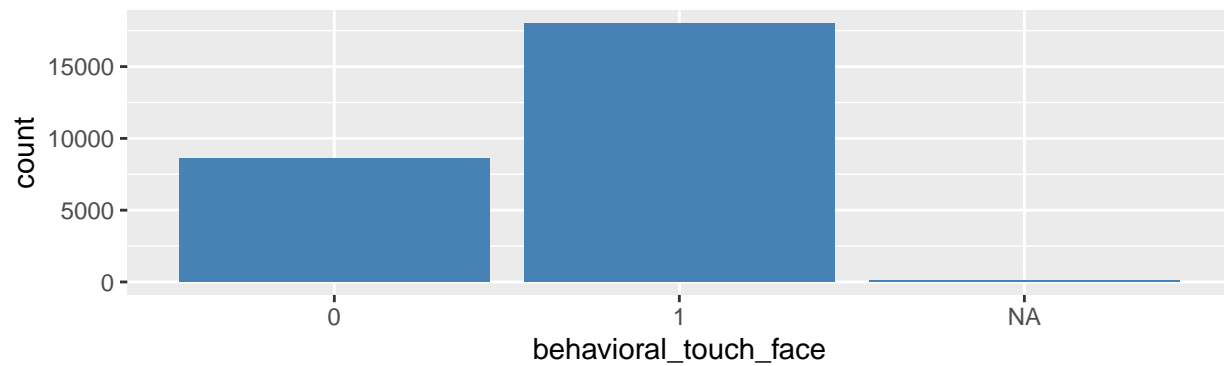
```
plots = list()
i = 1

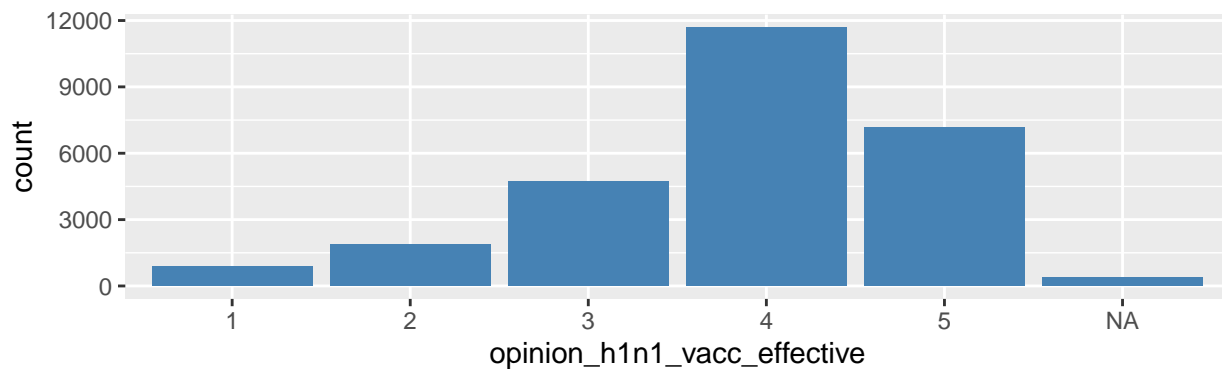
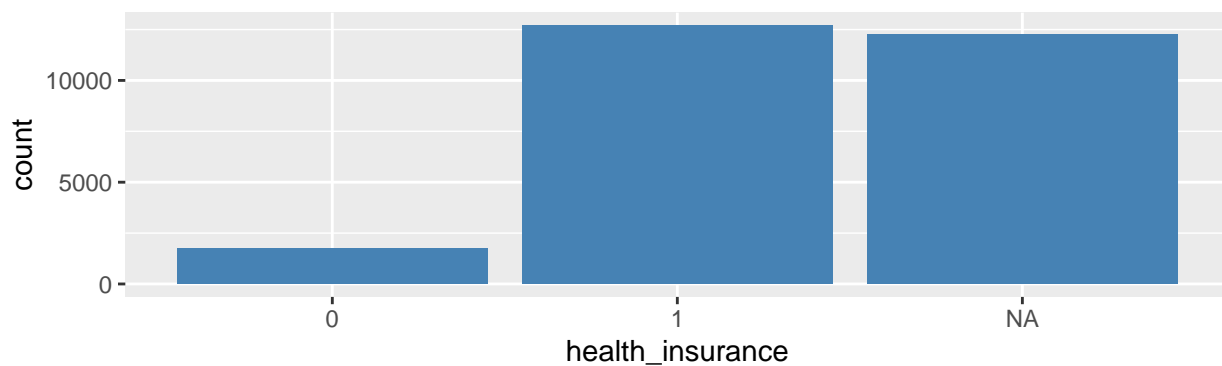
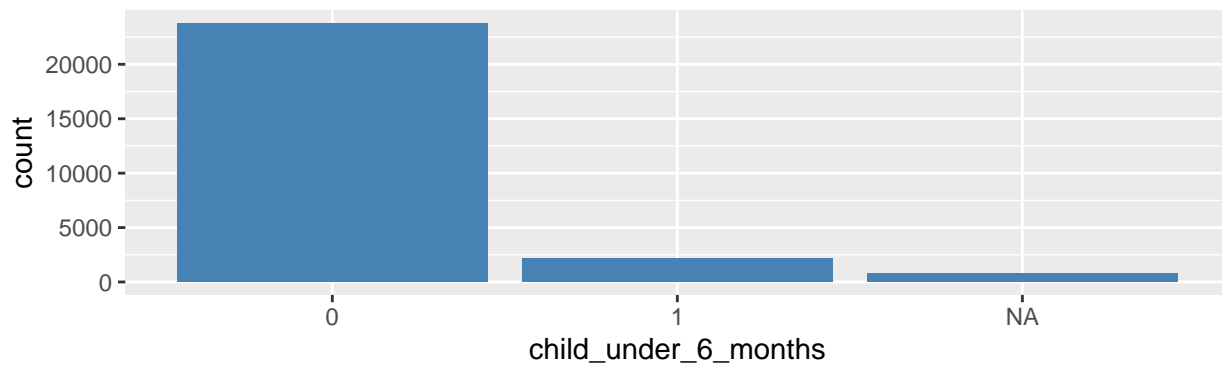
for(feature in features){
  p = df_train %>%
    ggplot(aes_string(x=feature)) +
    geom_bar(fill="steelblue")

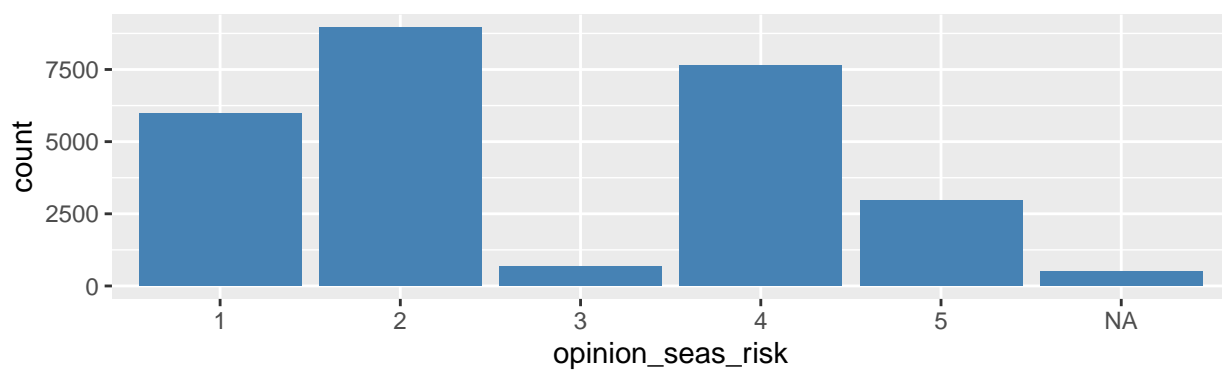
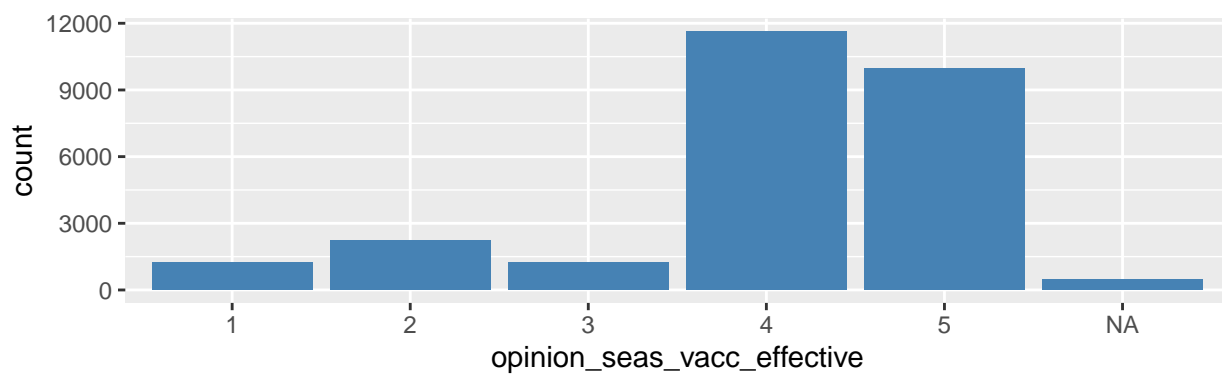
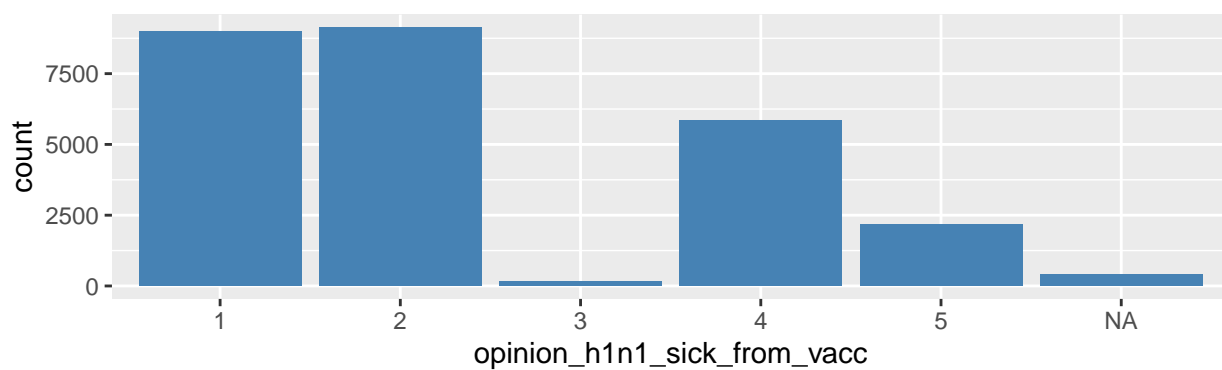
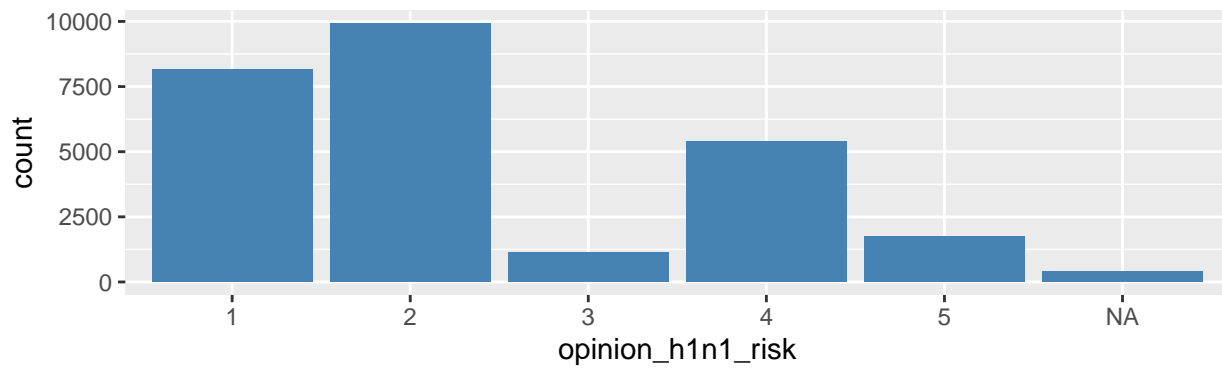
  plots[[i]] = p
  if(i %% 4 == 0){
    do.call("grid.arrange", c(plots[(i-3): i], nrow=4))
  }
  i = i + 1
}
```

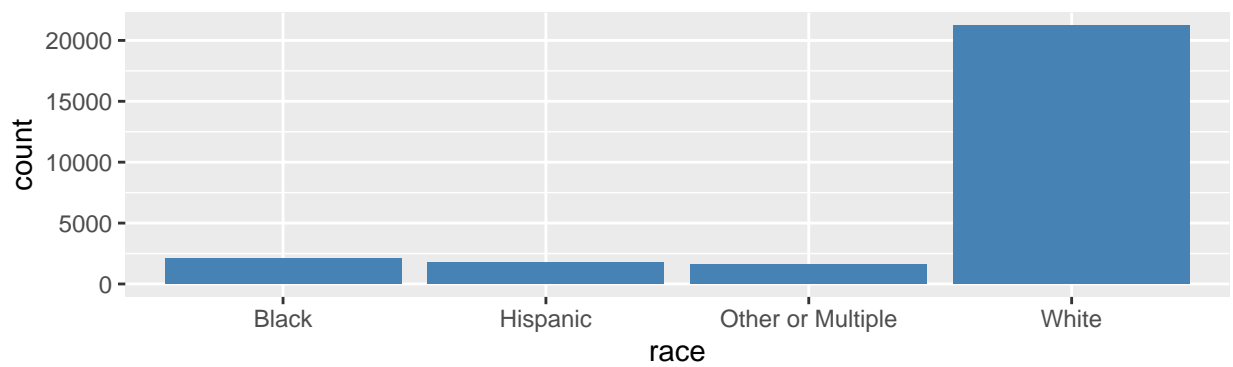
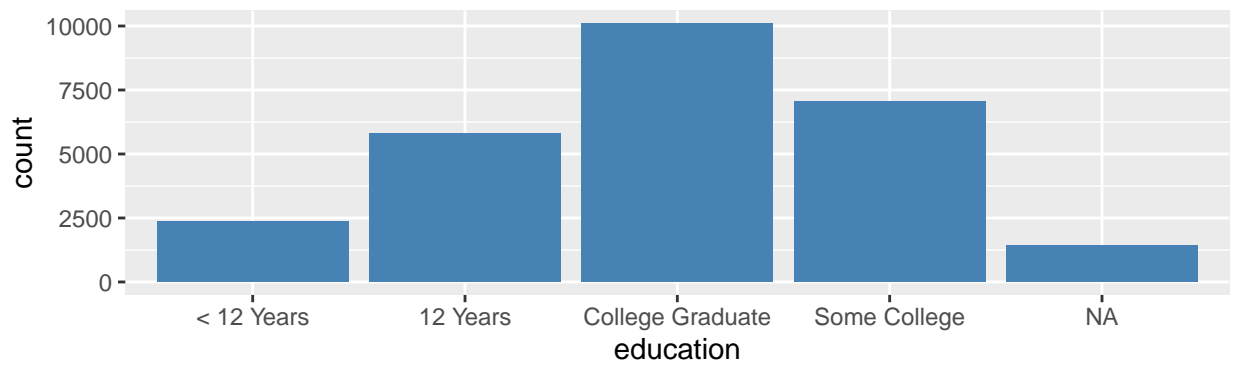
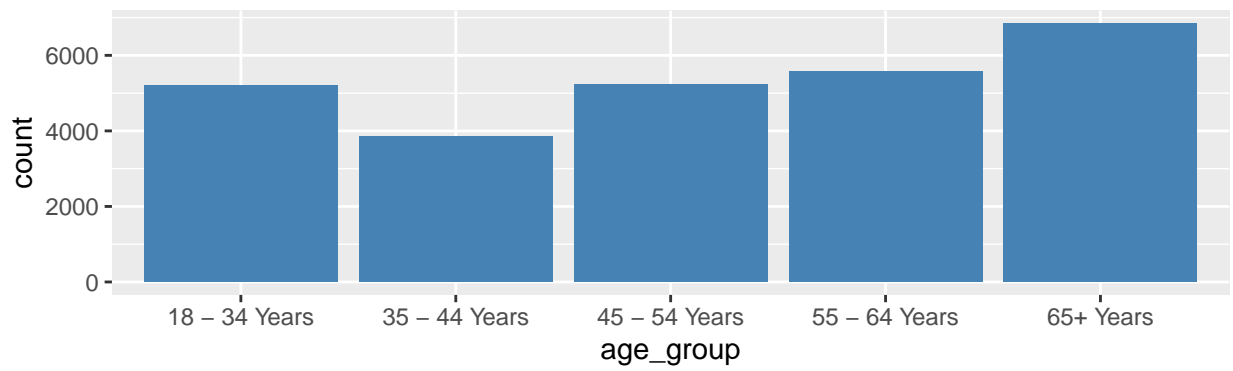
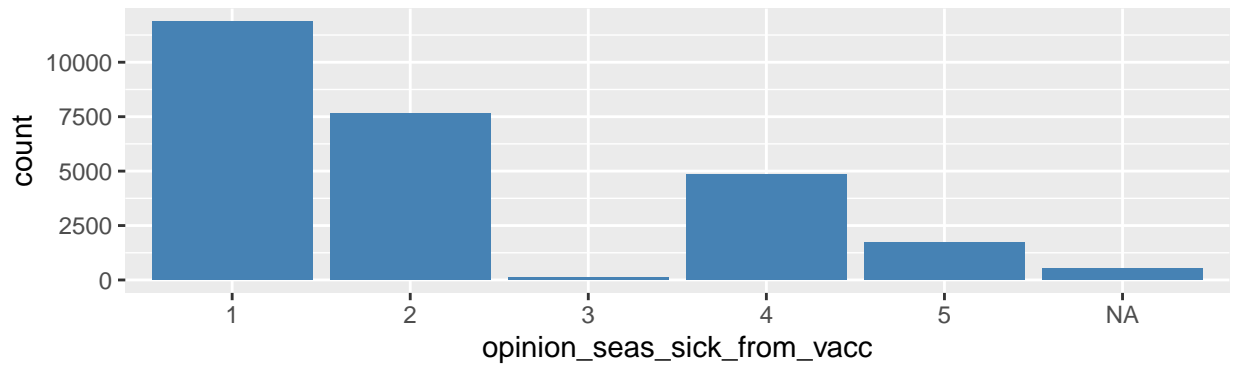


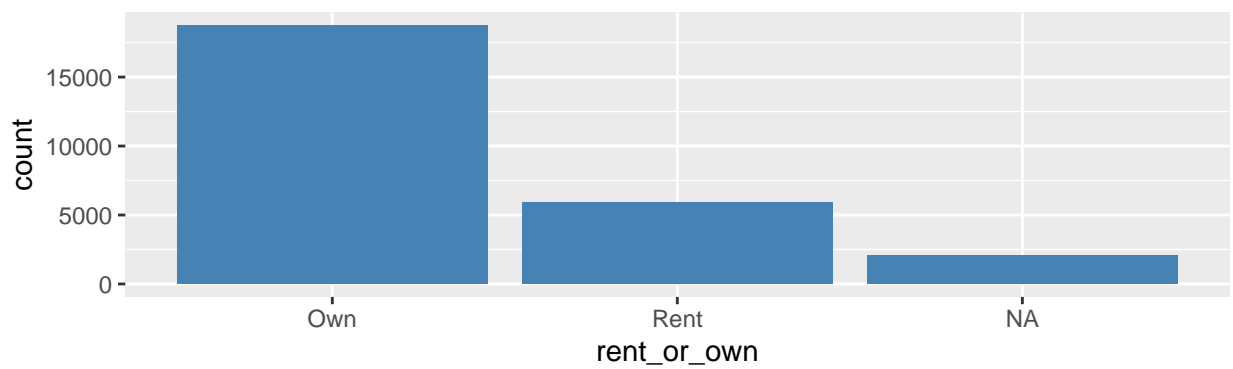
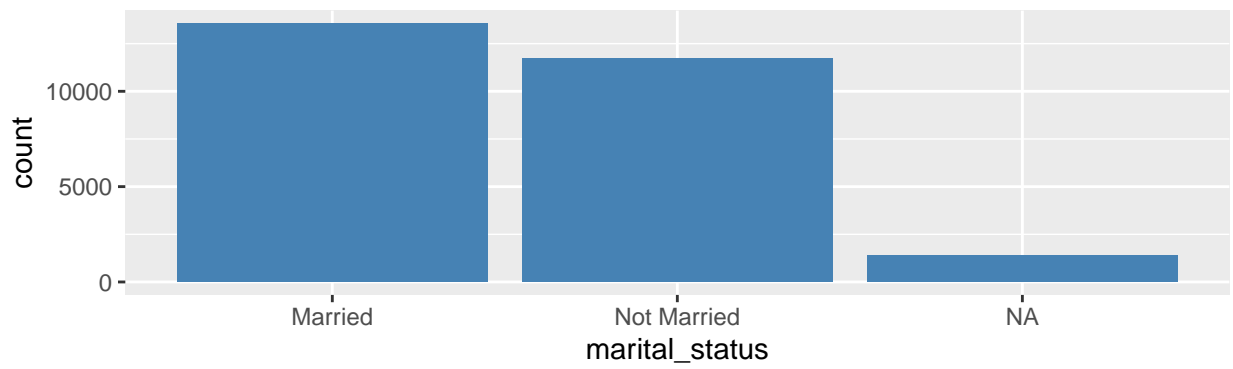
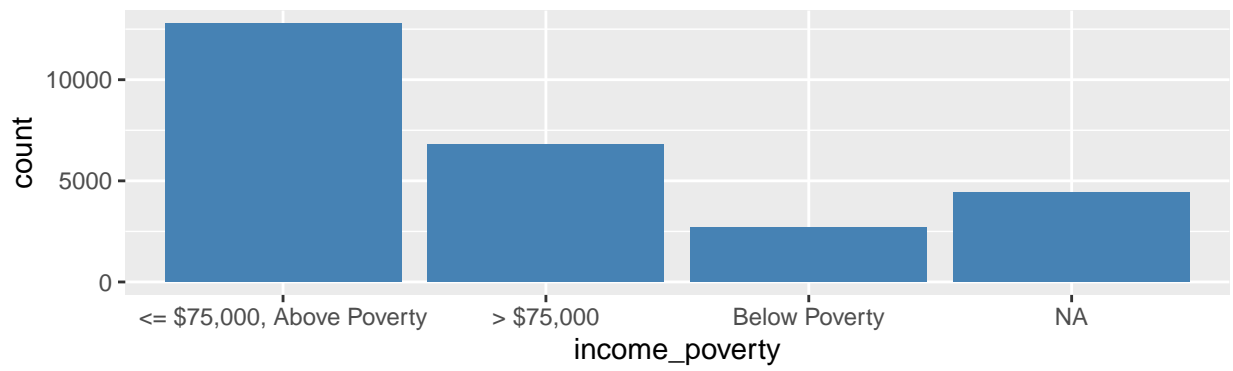
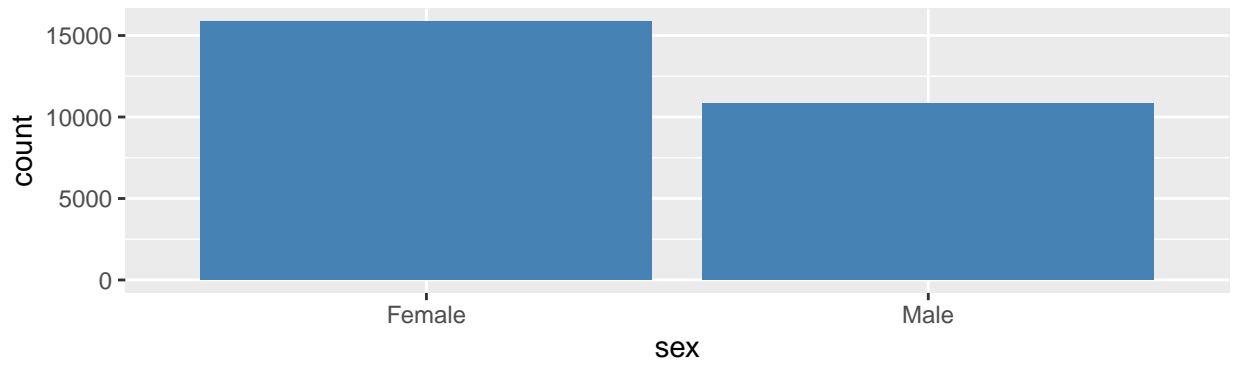


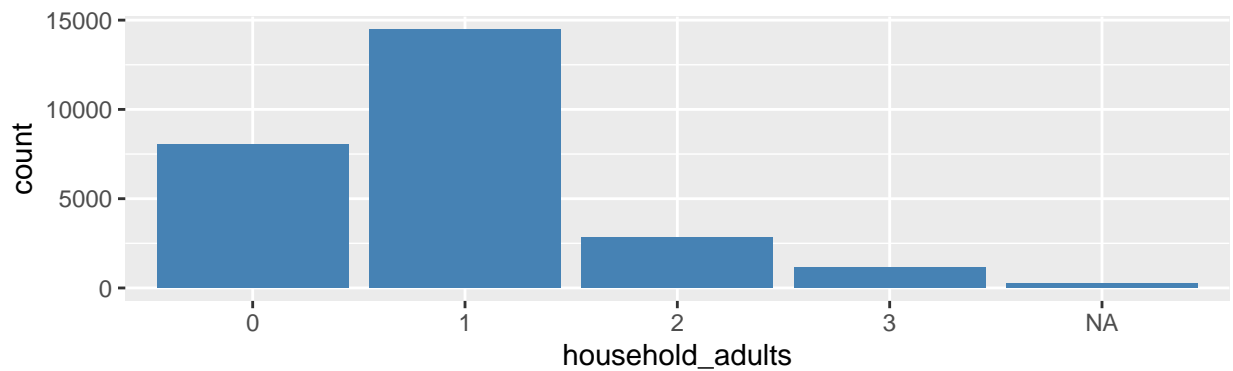
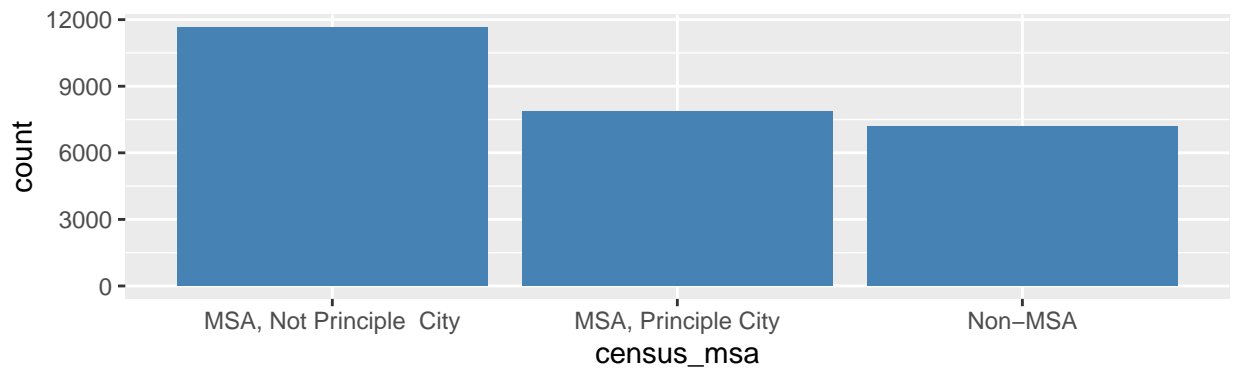
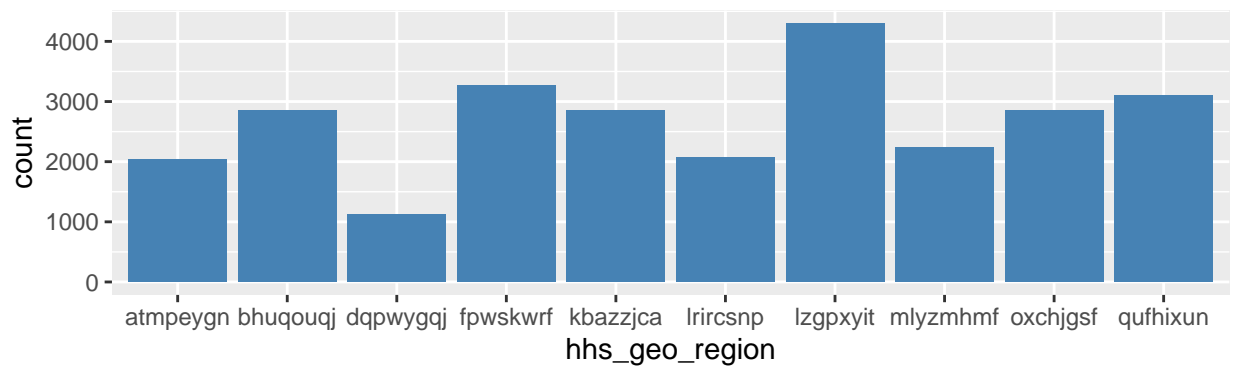
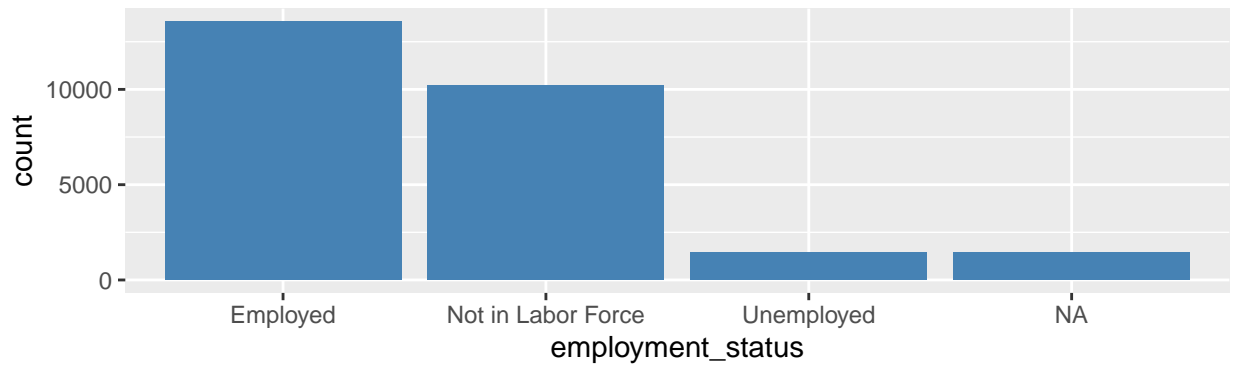








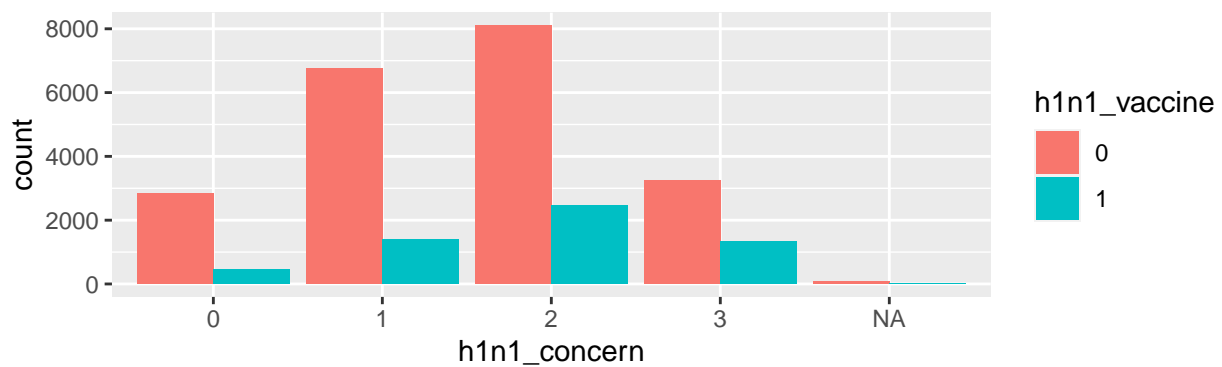
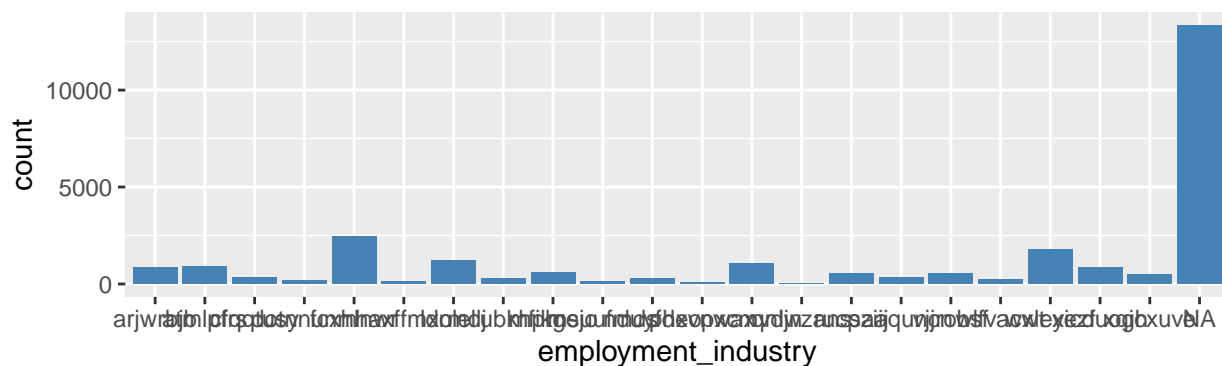
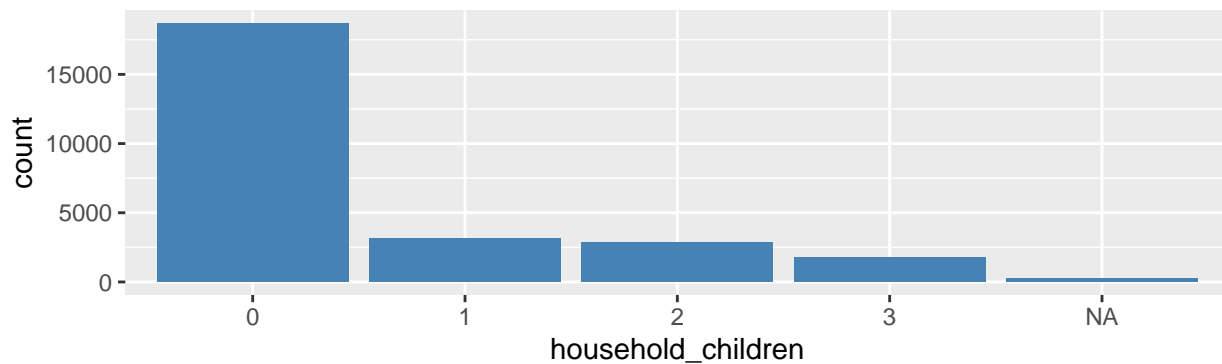


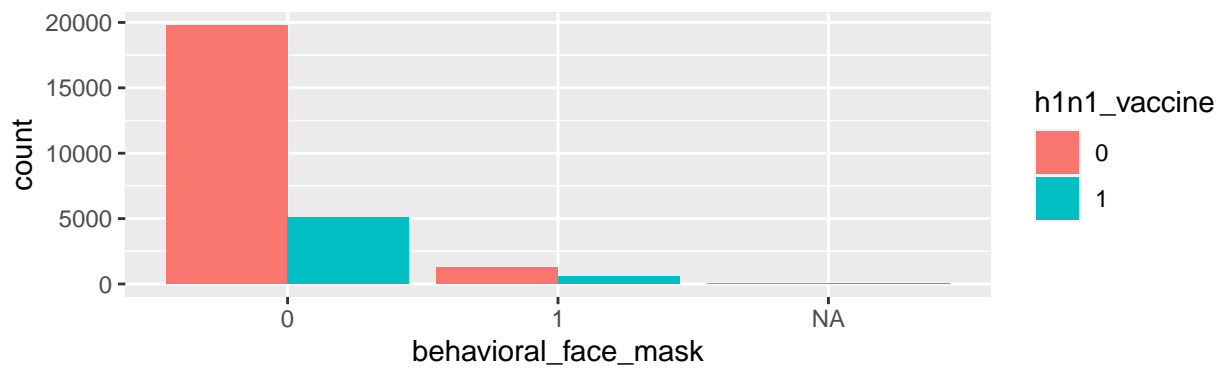
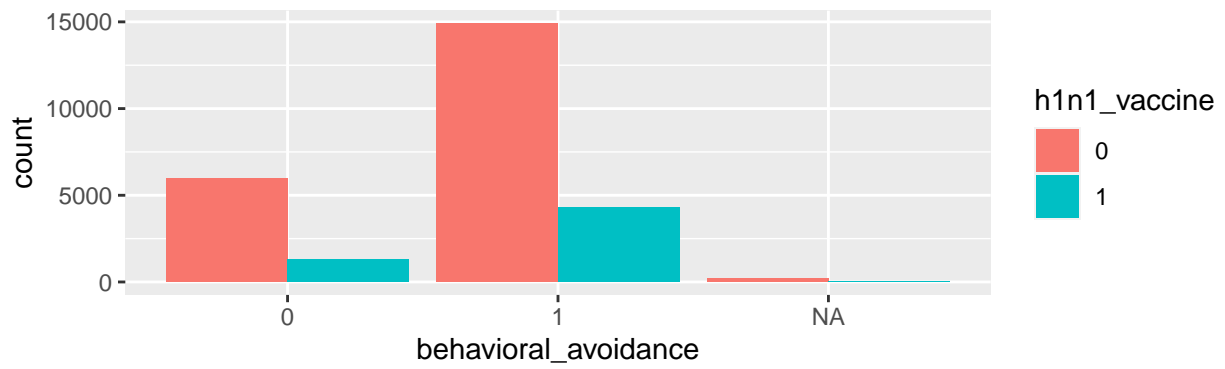
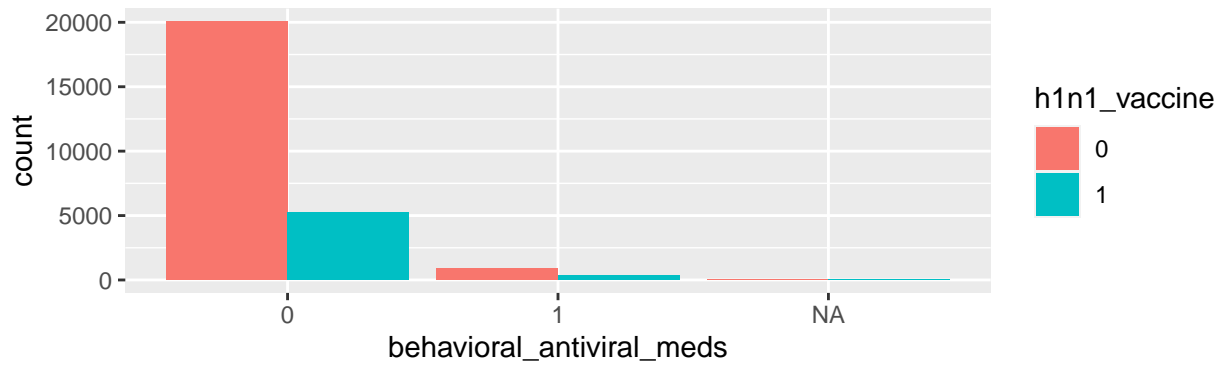
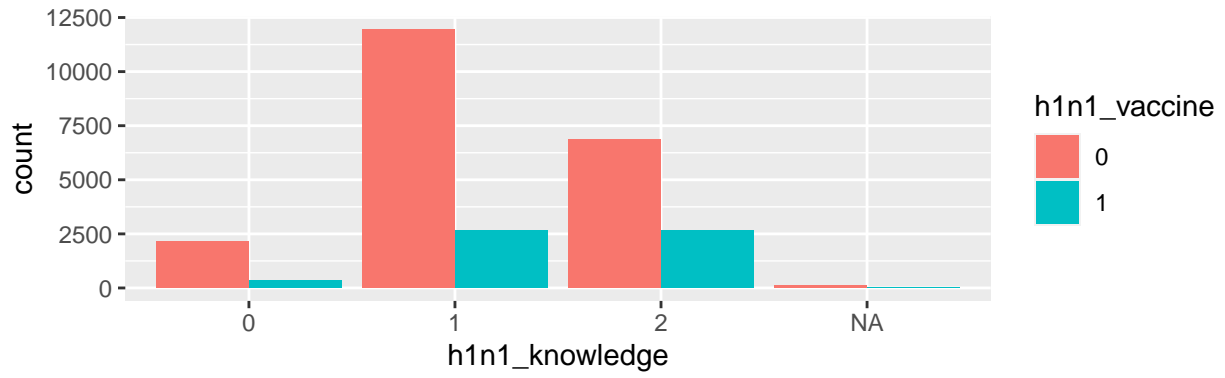


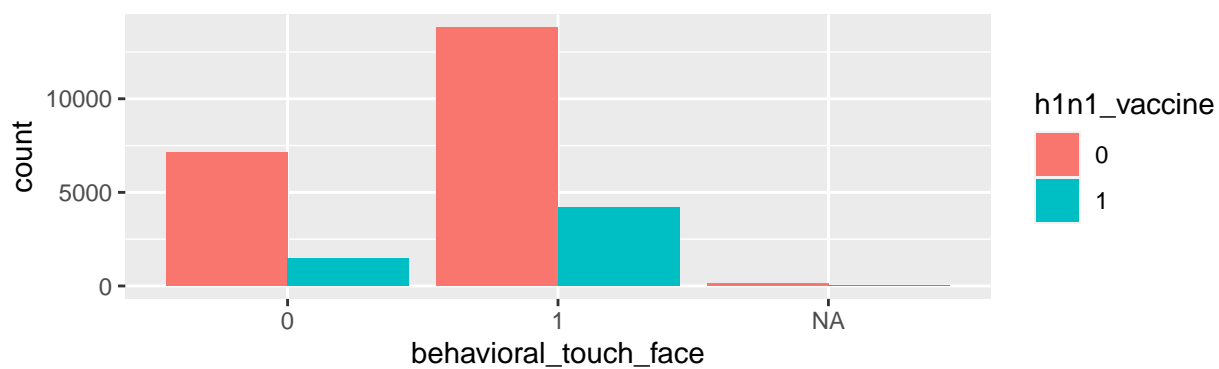
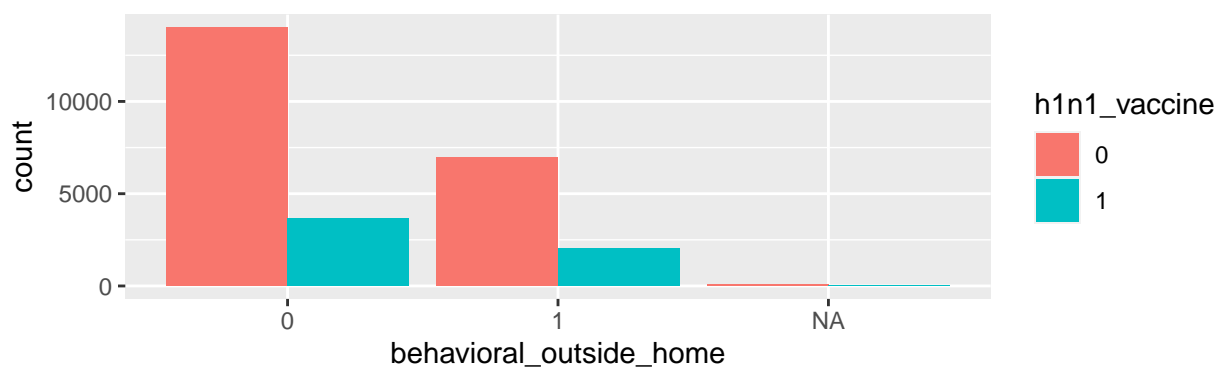
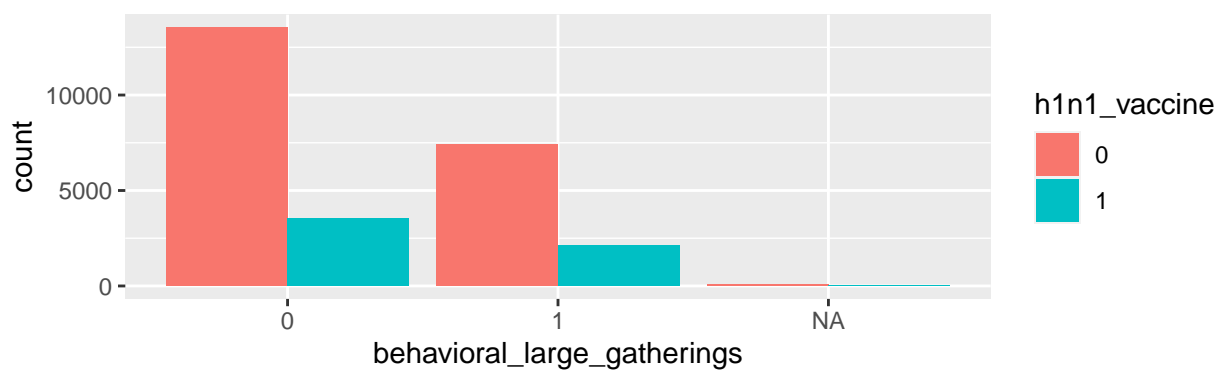
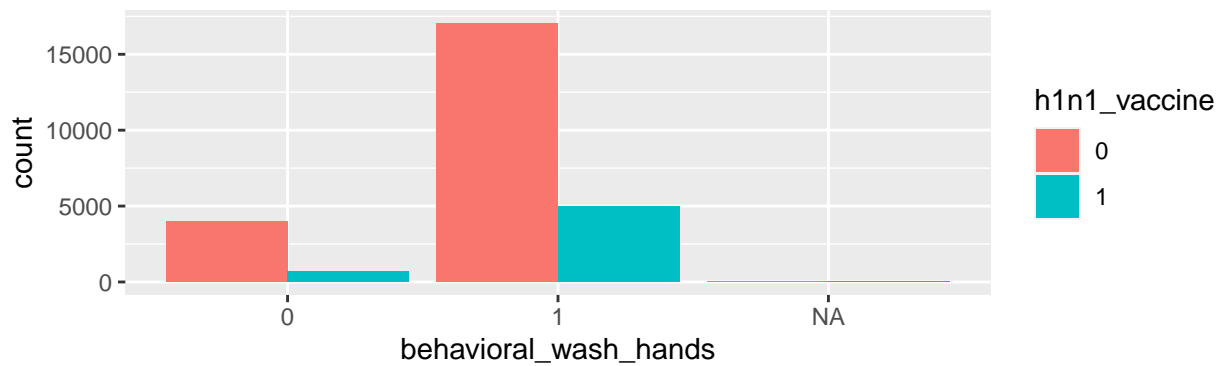
Distribution of each feature by target

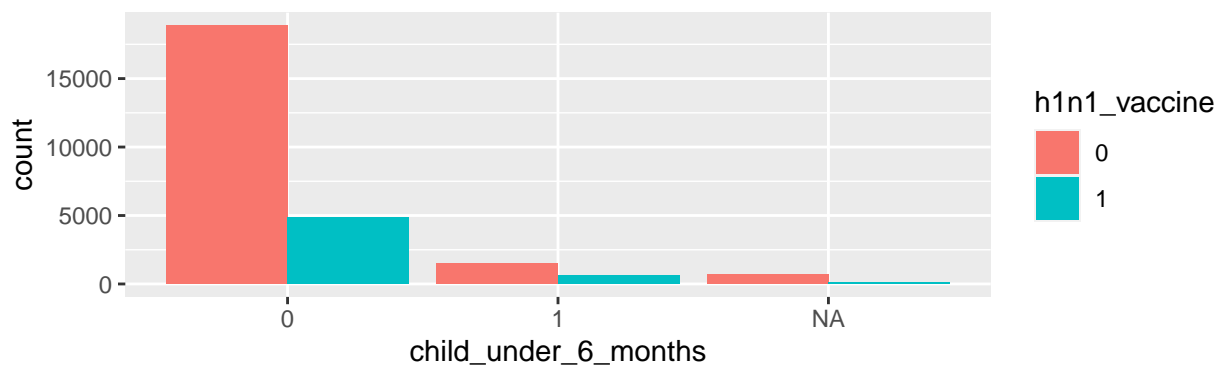
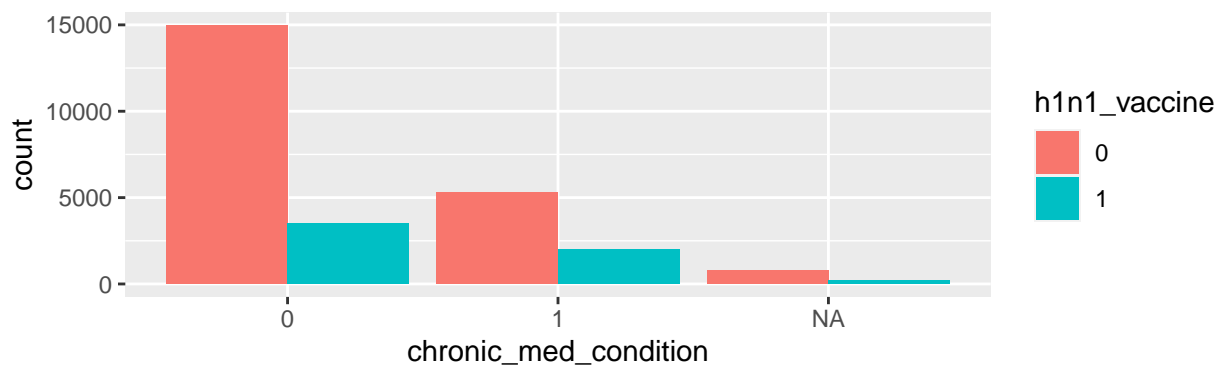
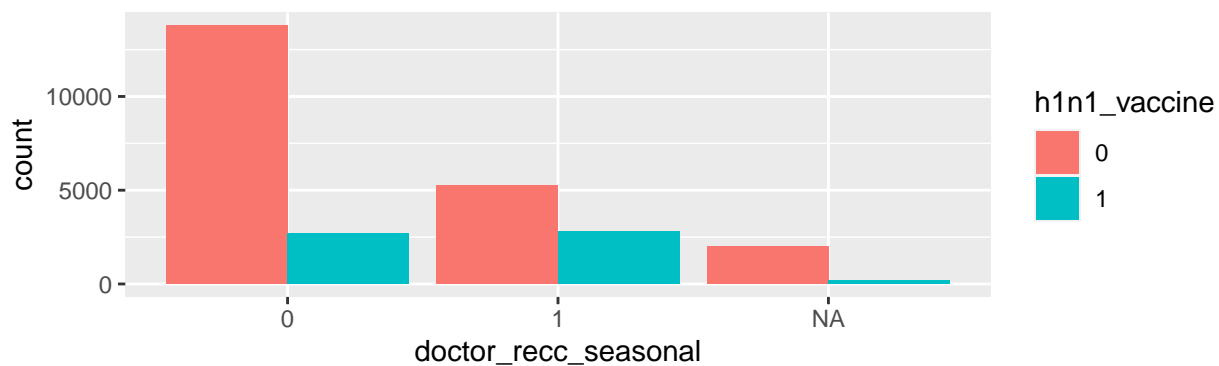
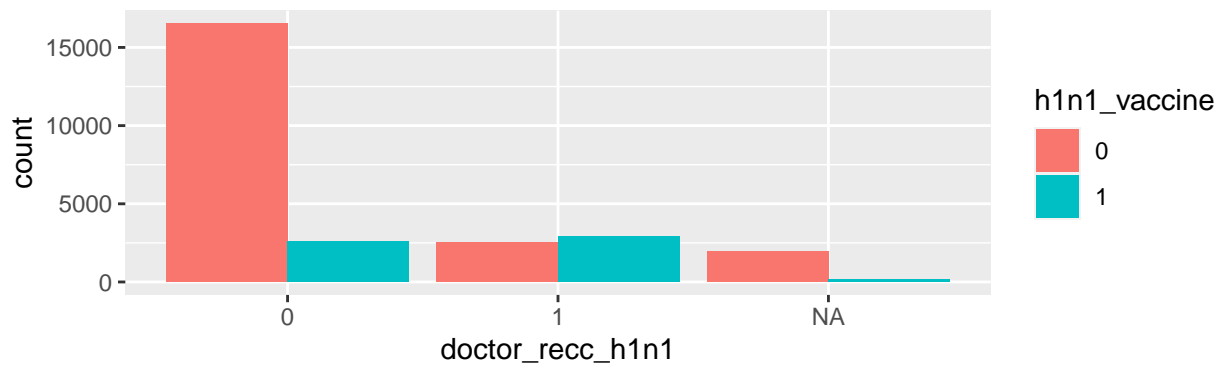
```
for(target in targets){
  for(feature in features){
    p = df_train %>%
```

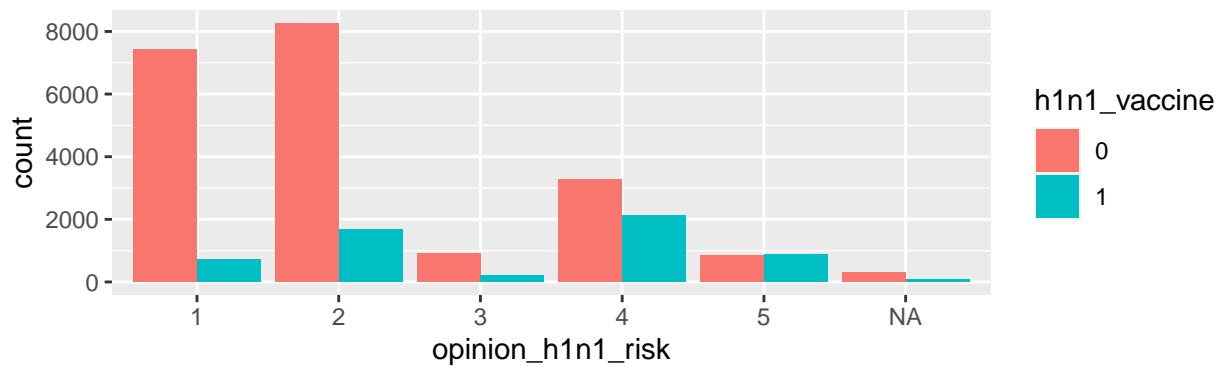
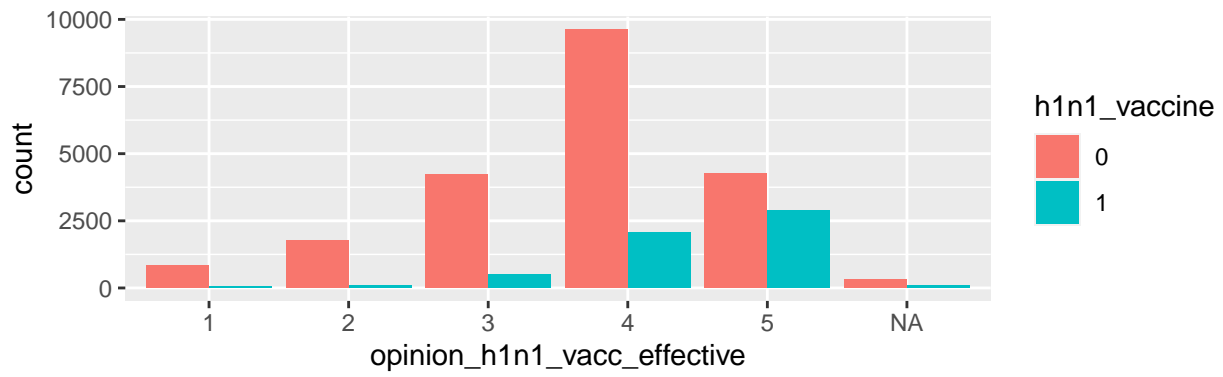
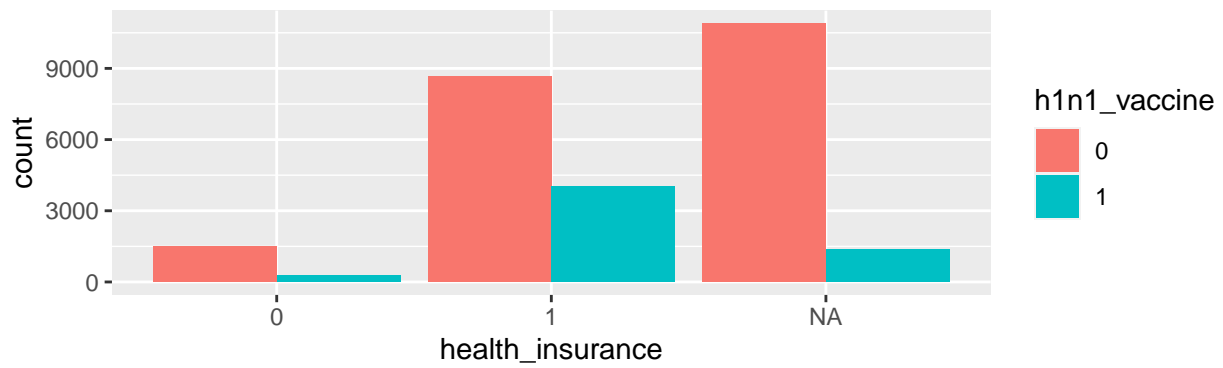
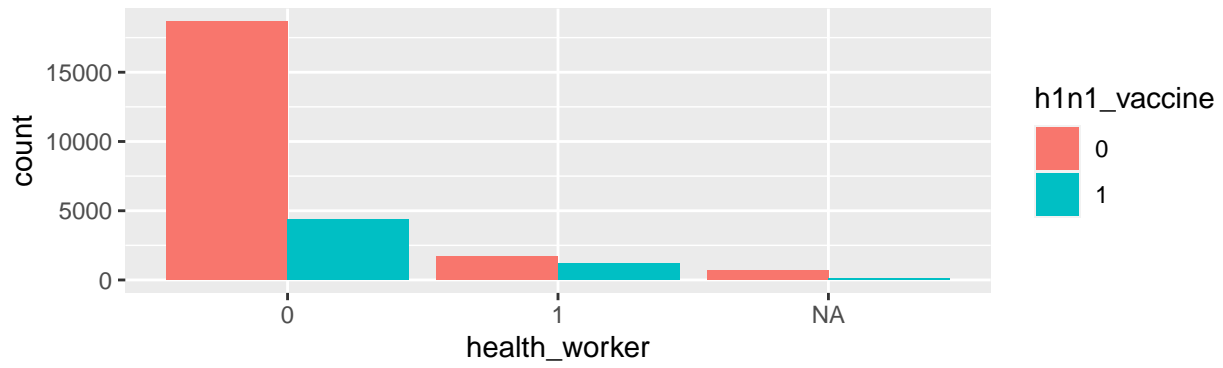
```
ggplot(aes_string(x=feature, fill=target)) +  
  geom_bar(position="dodge")  
  
plots[[i]] = p  
if(i %% 4 == 0){  
  do.call("grid.arrange", c(plots[(i-3): i], nrow=4))  
}  
i = i + 1  
}  
}
```

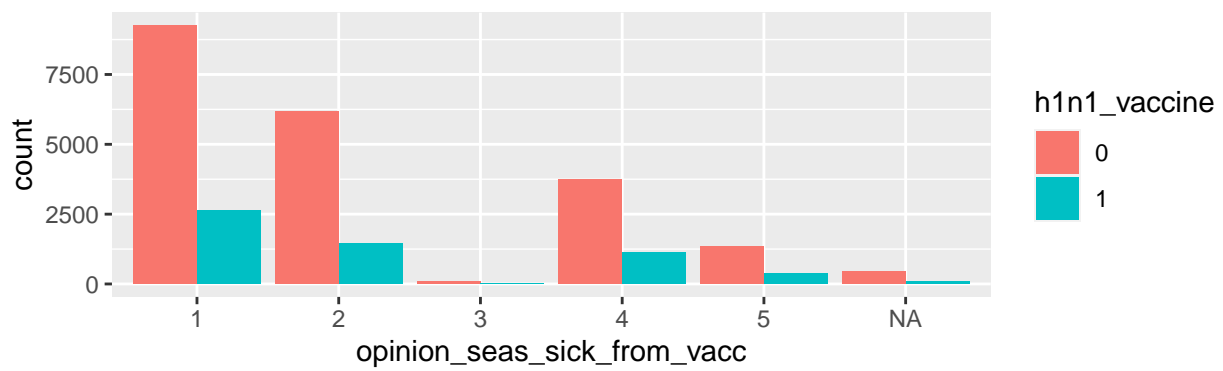
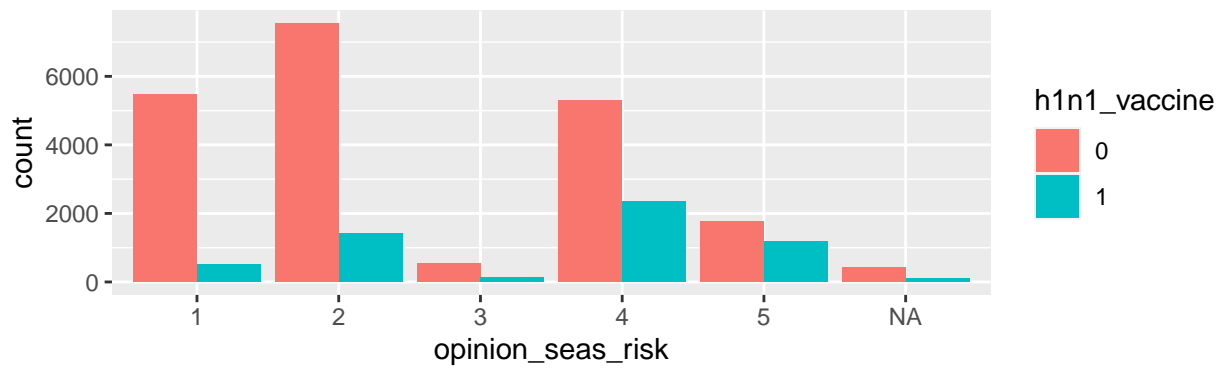
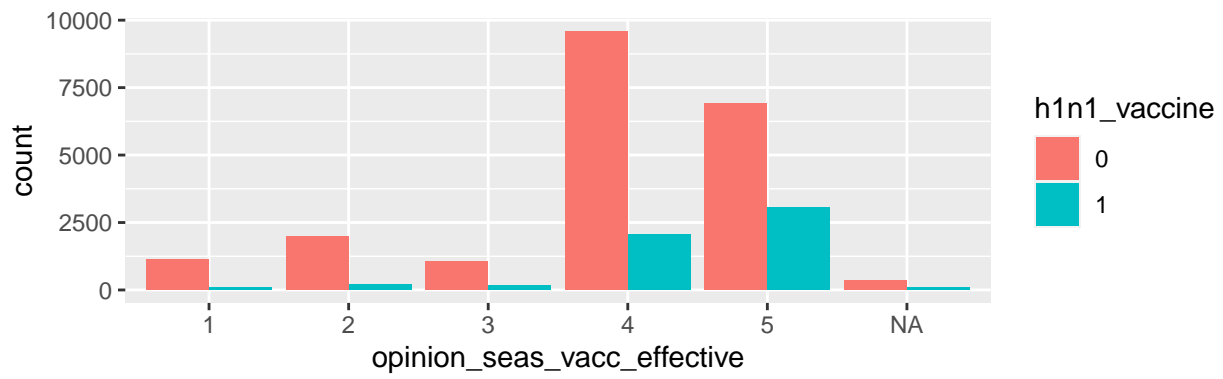
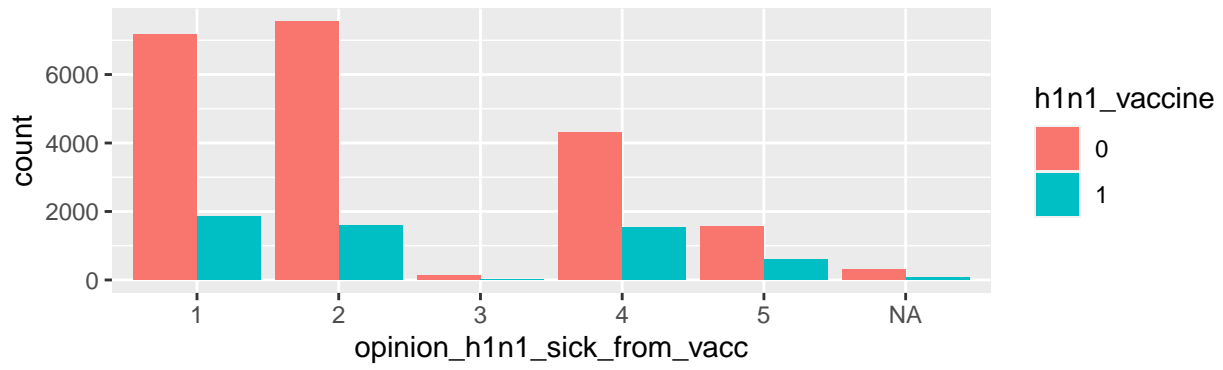


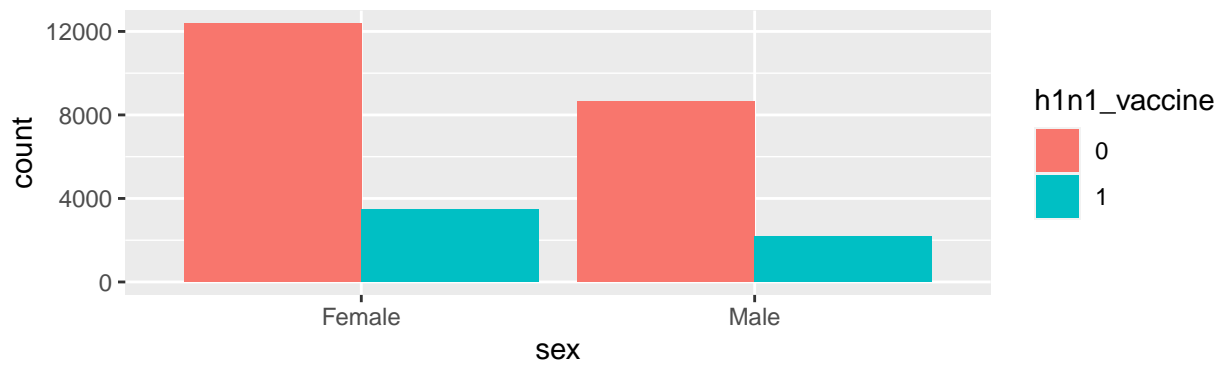
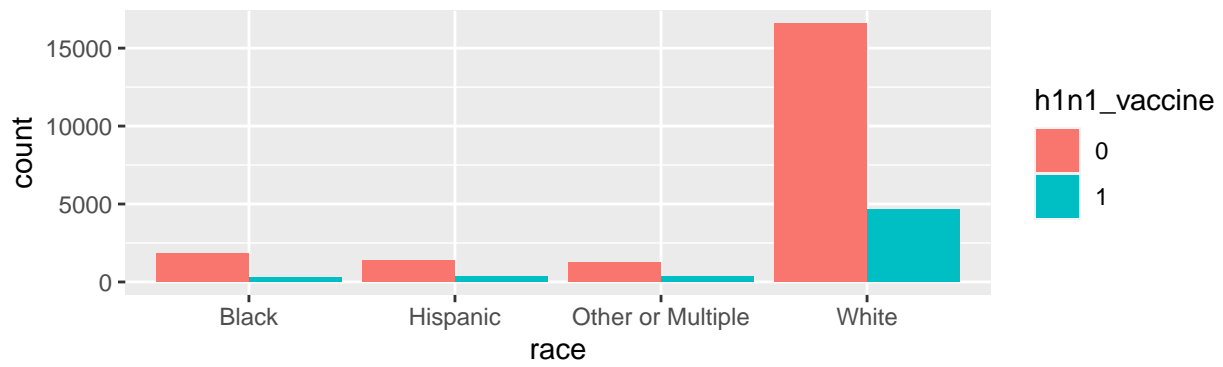
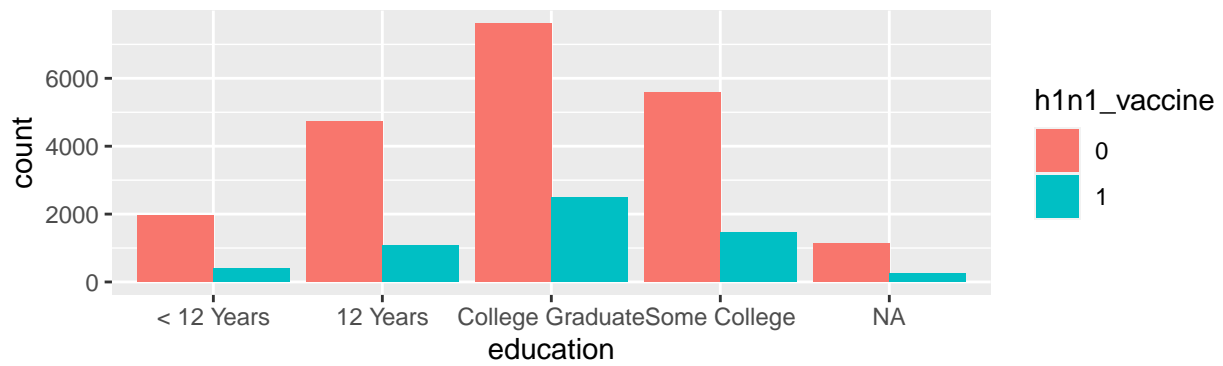
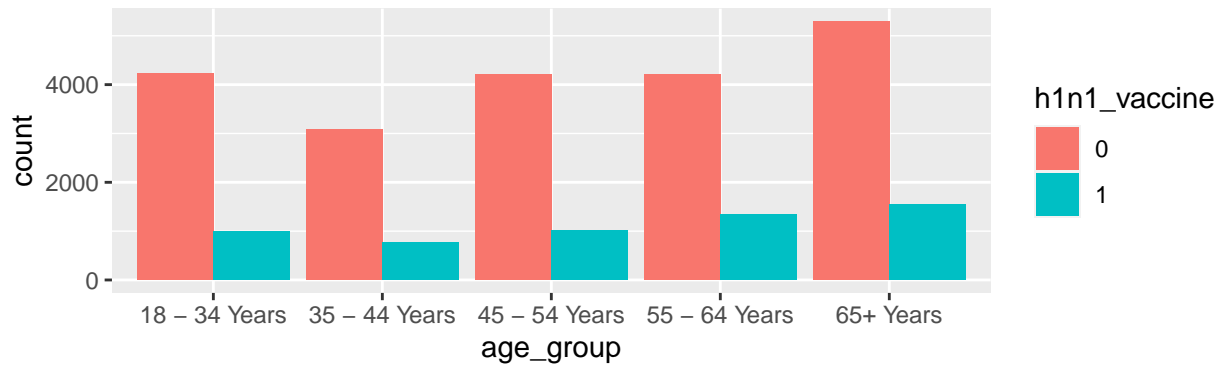


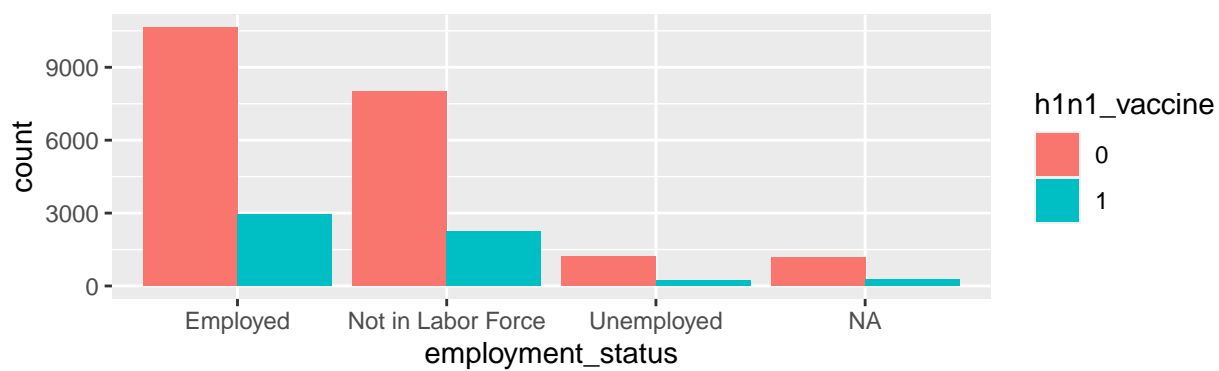
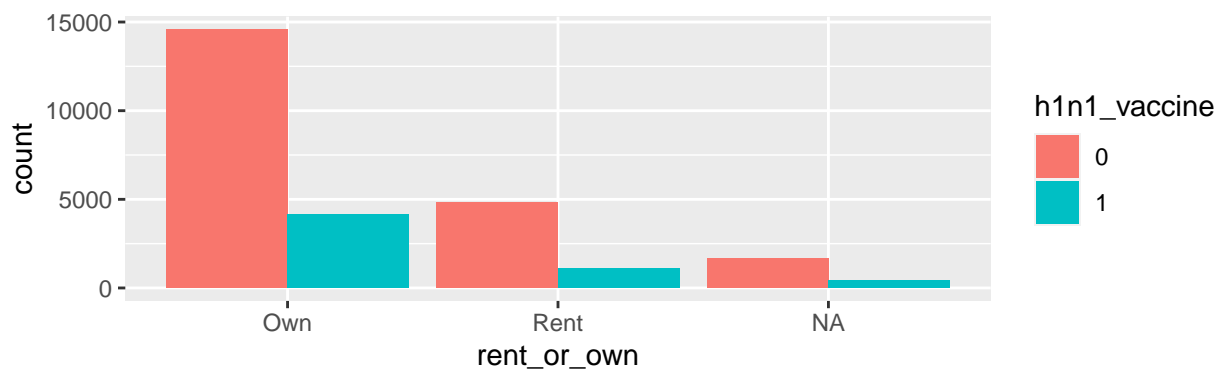
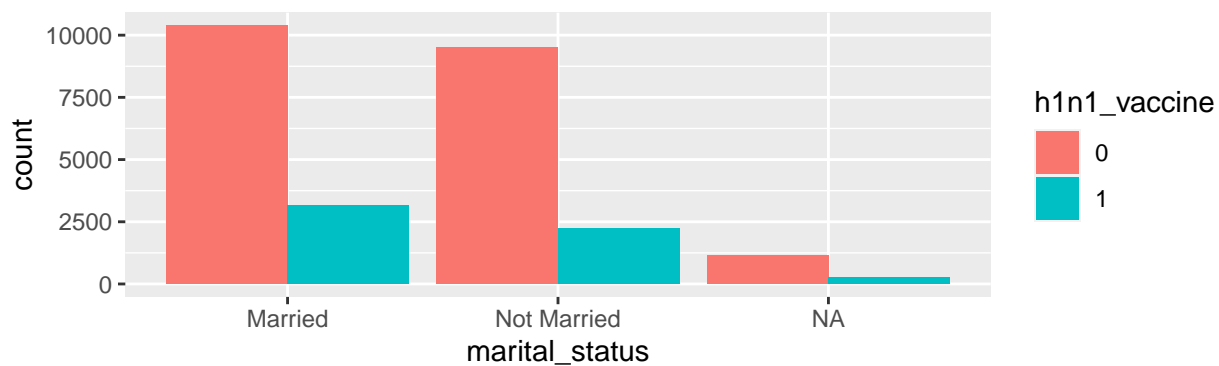
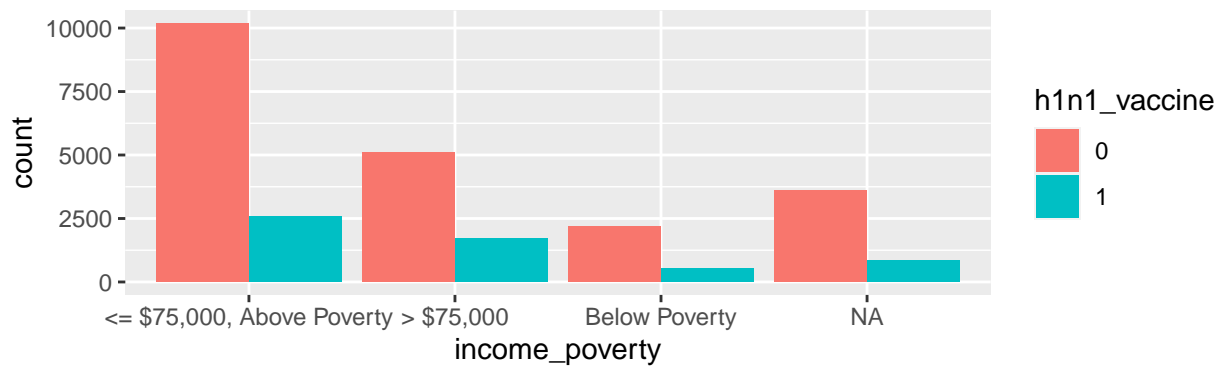


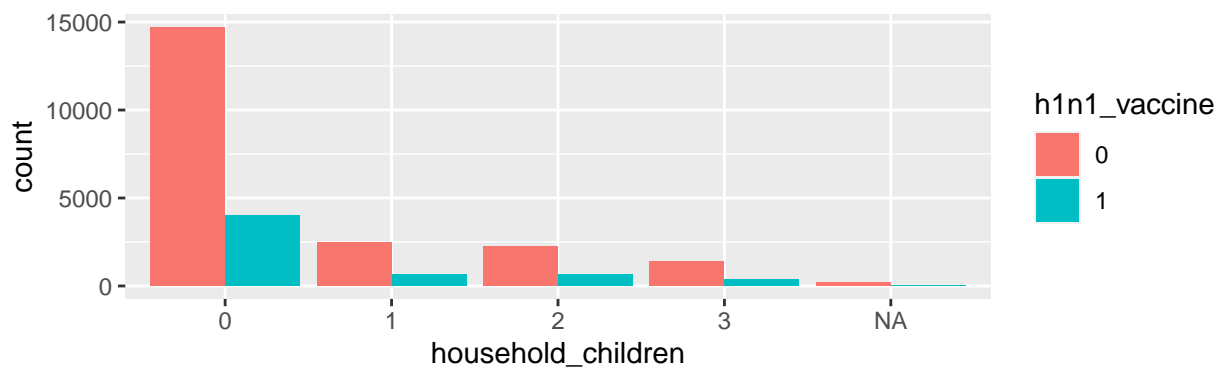
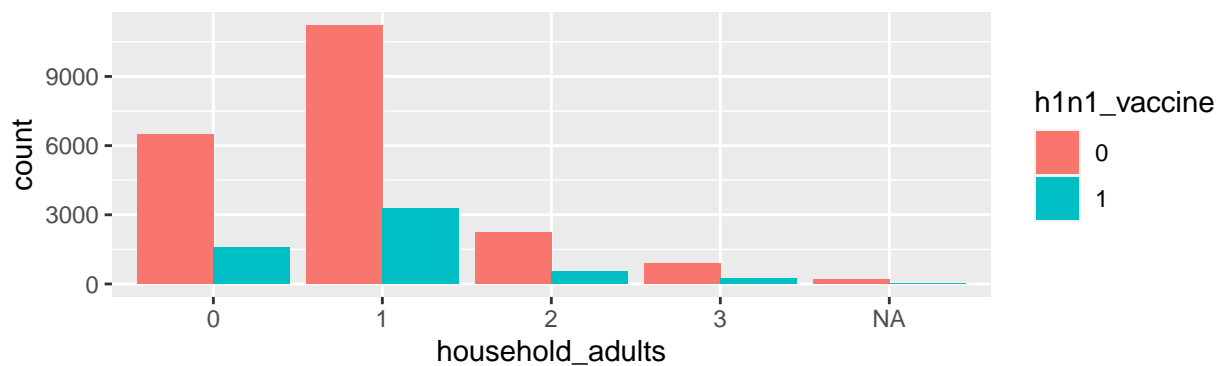
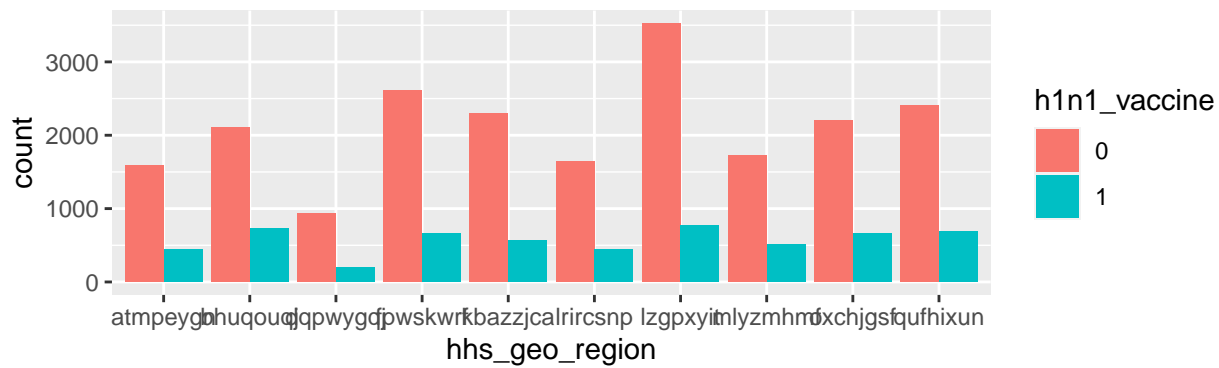


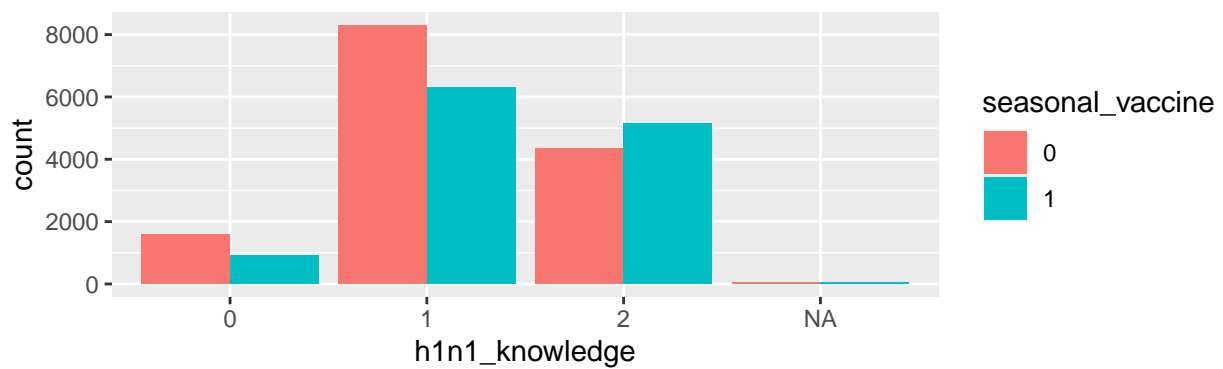
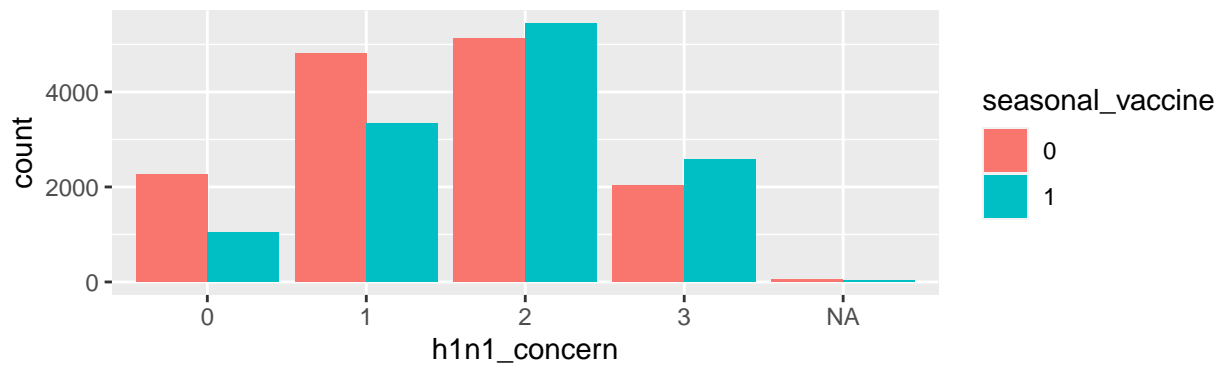
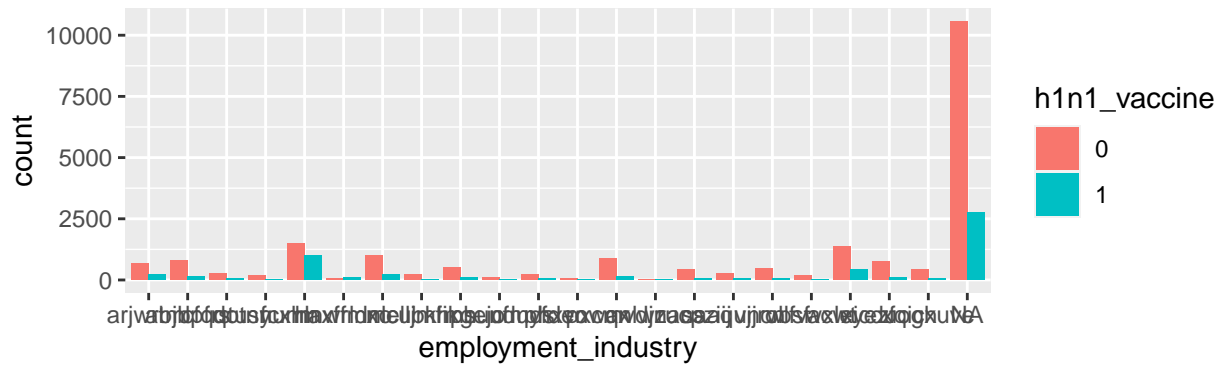


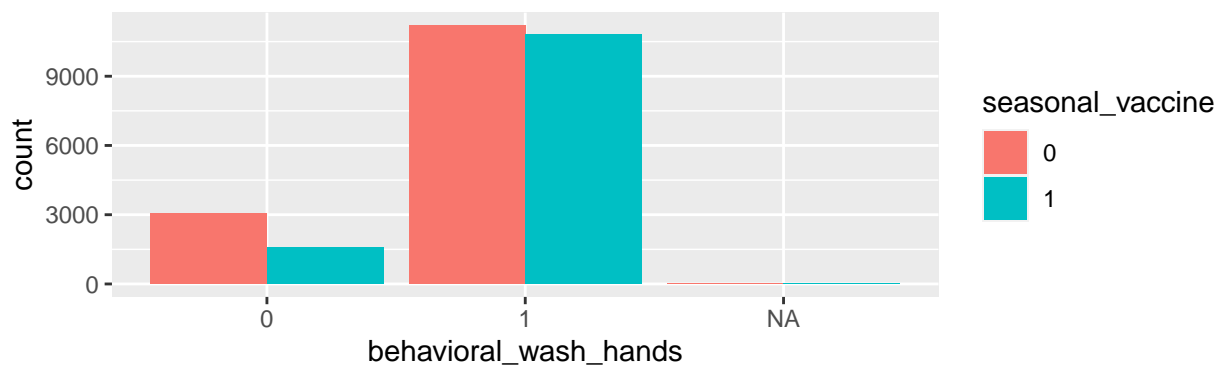
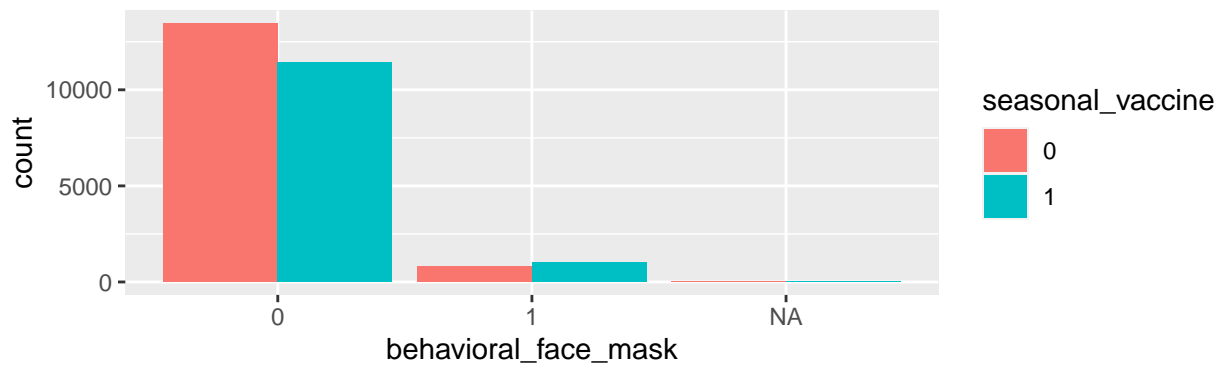
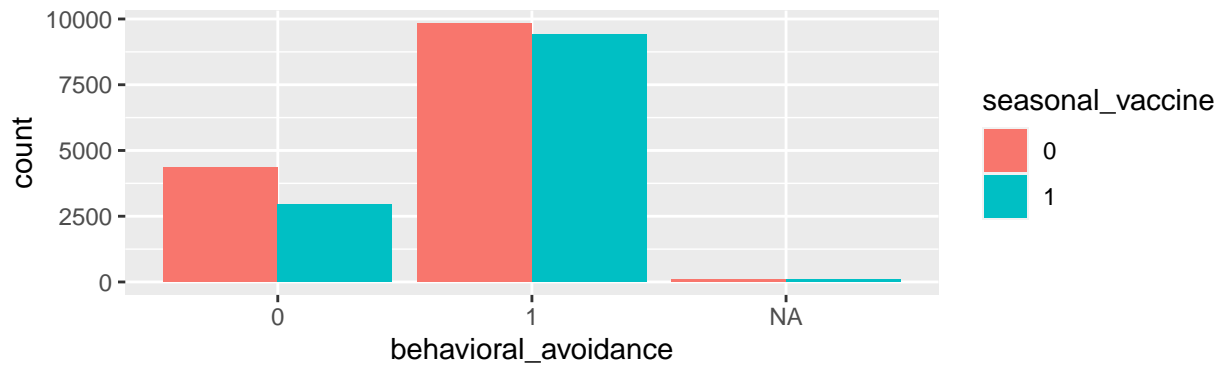
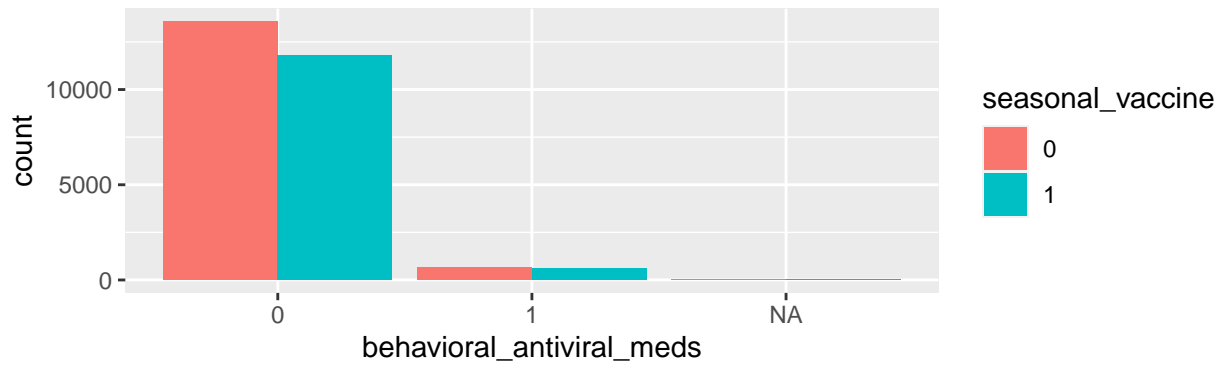


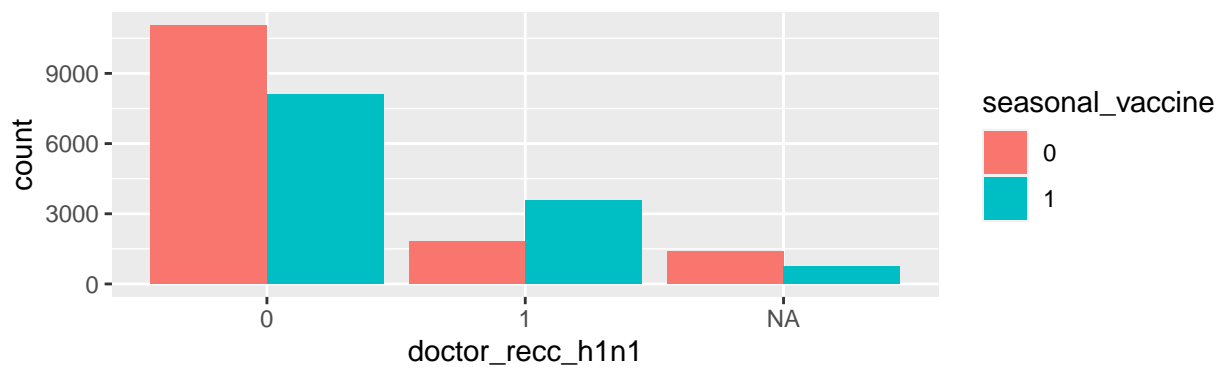
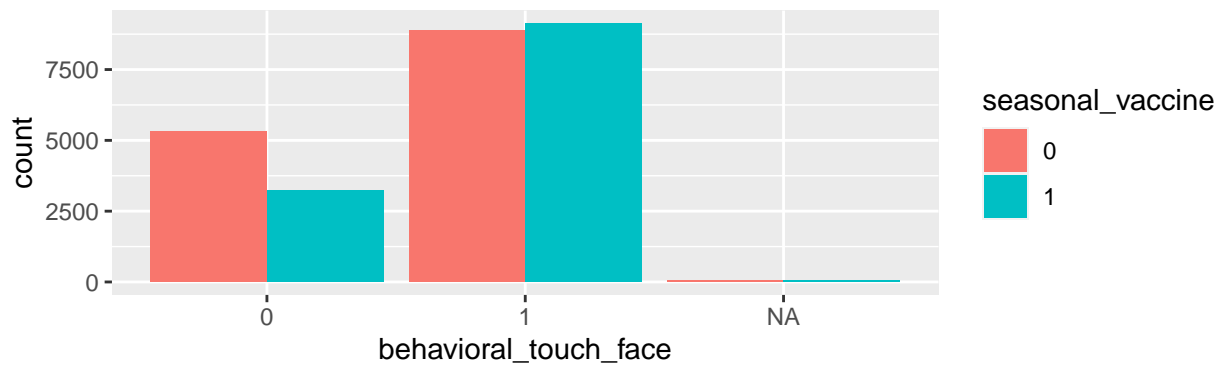
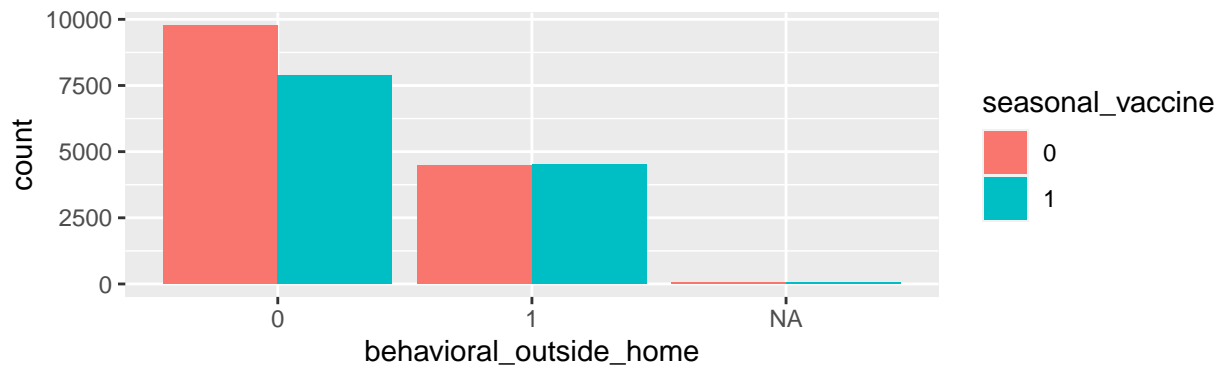
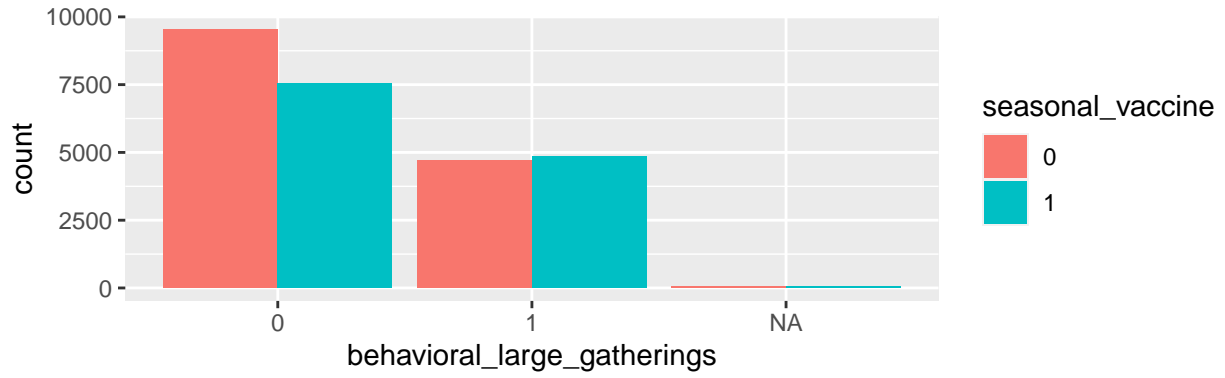


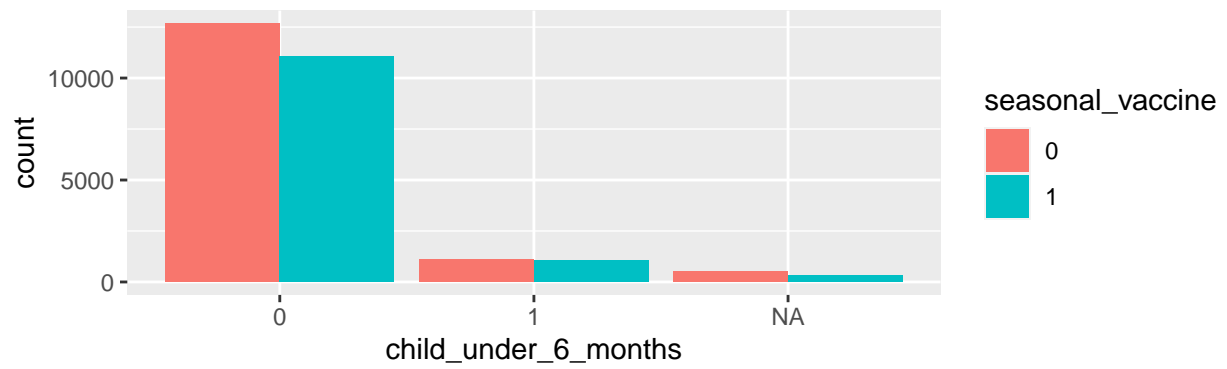
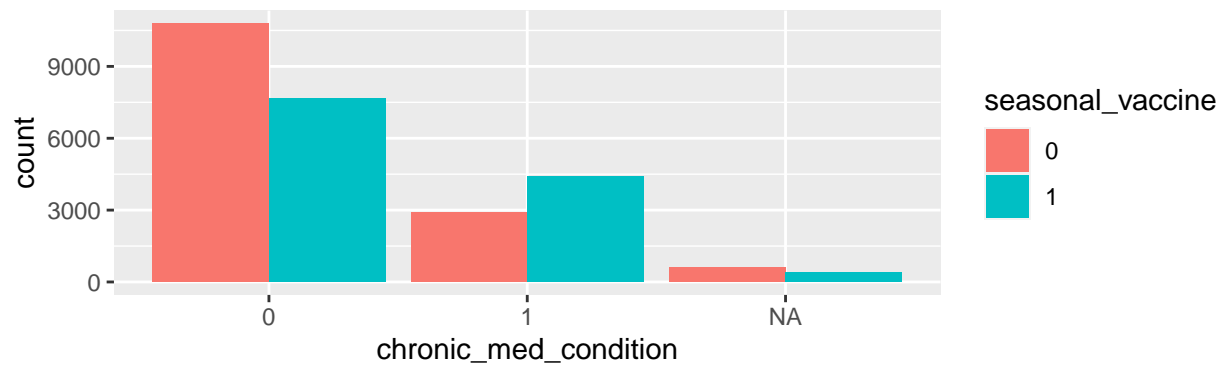
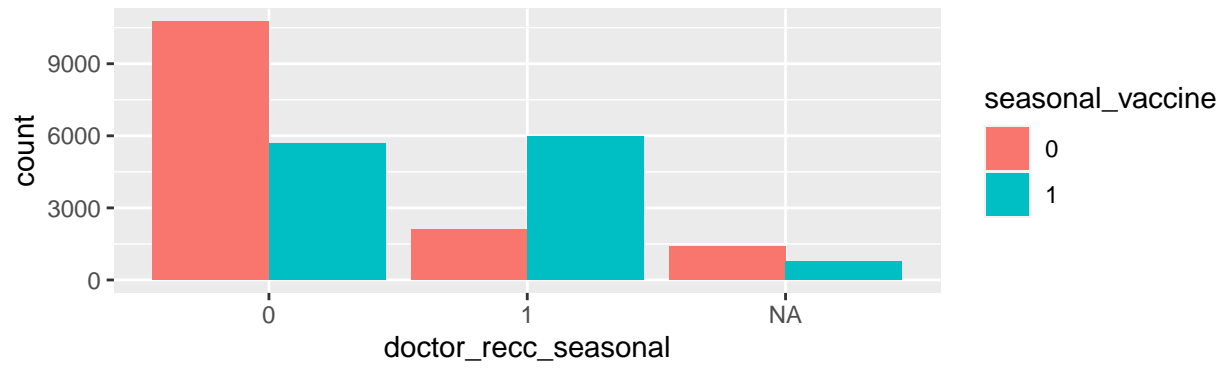


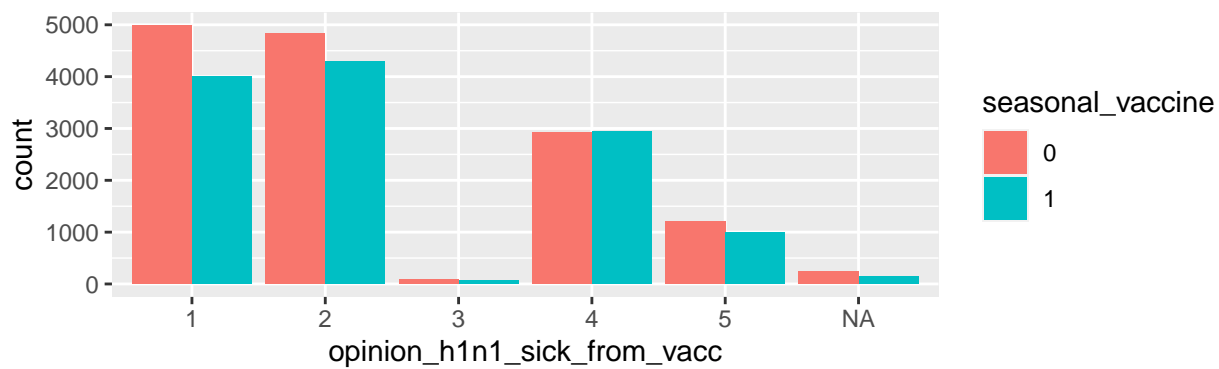
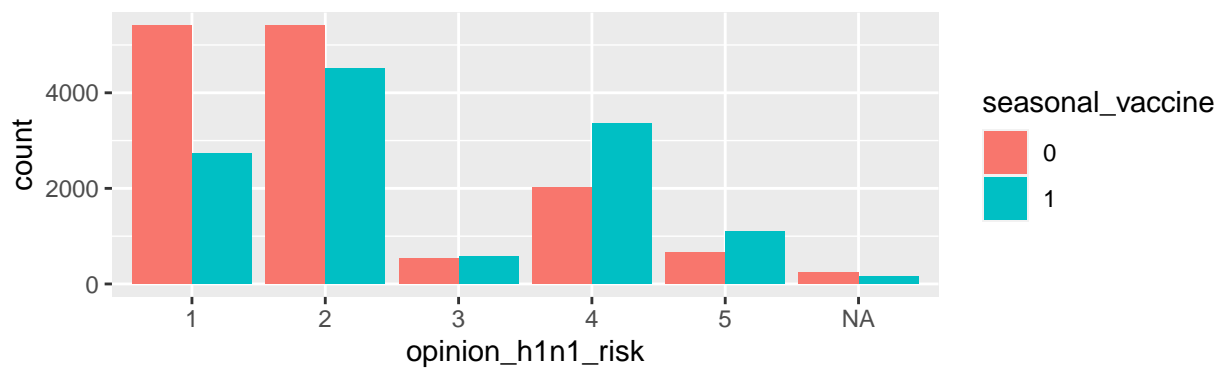
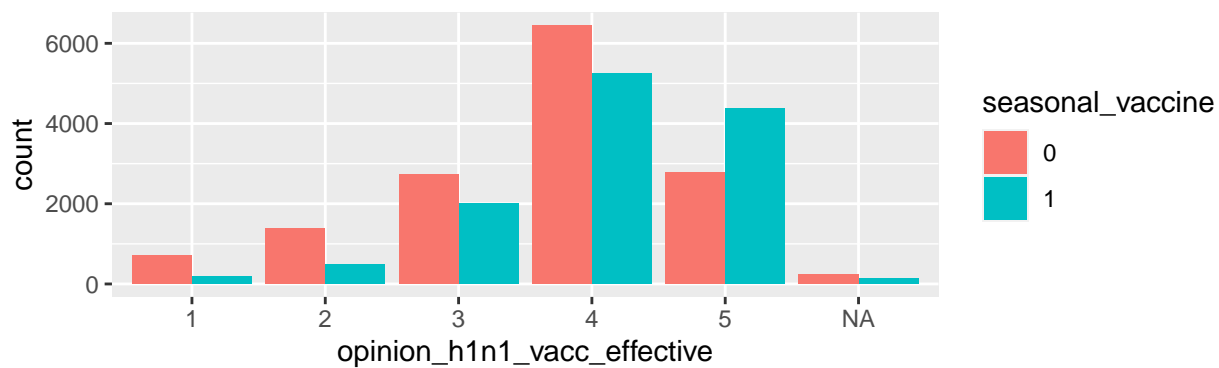
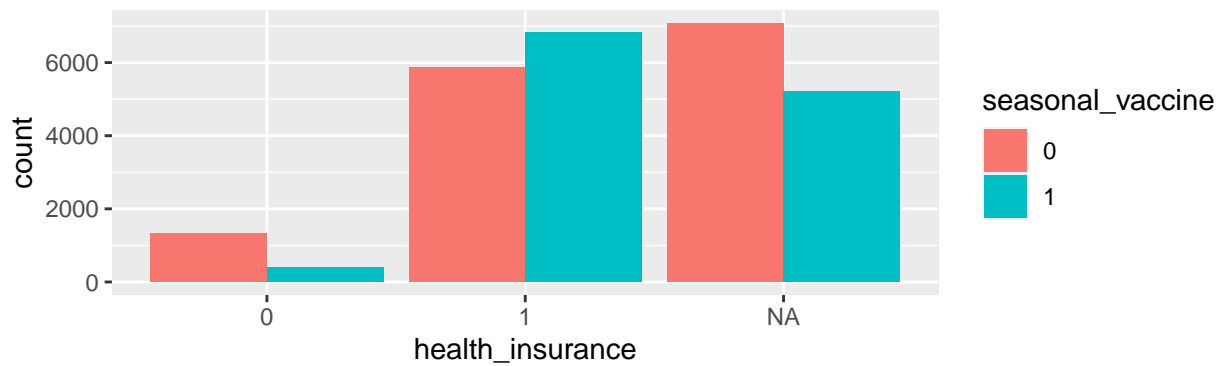


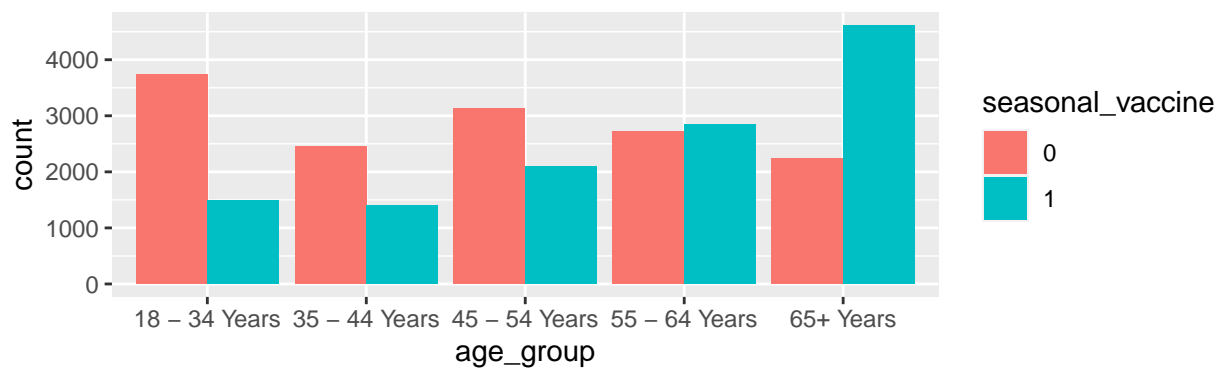
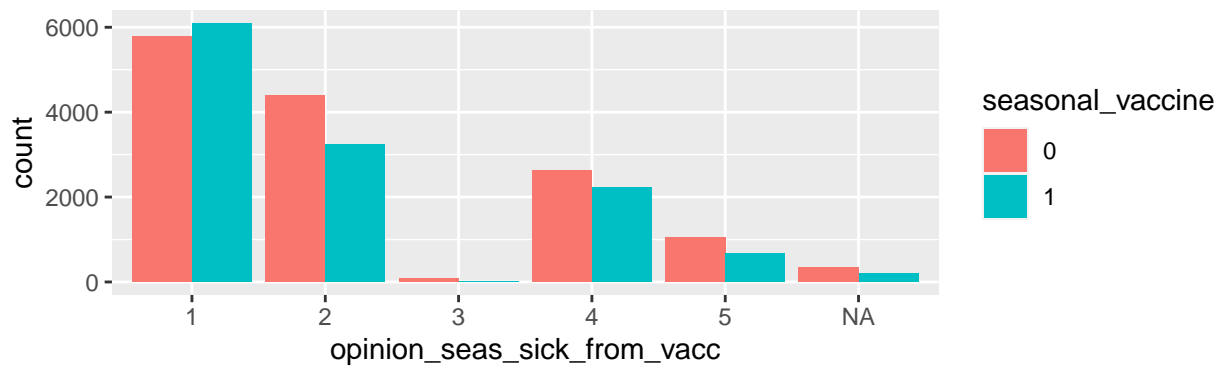
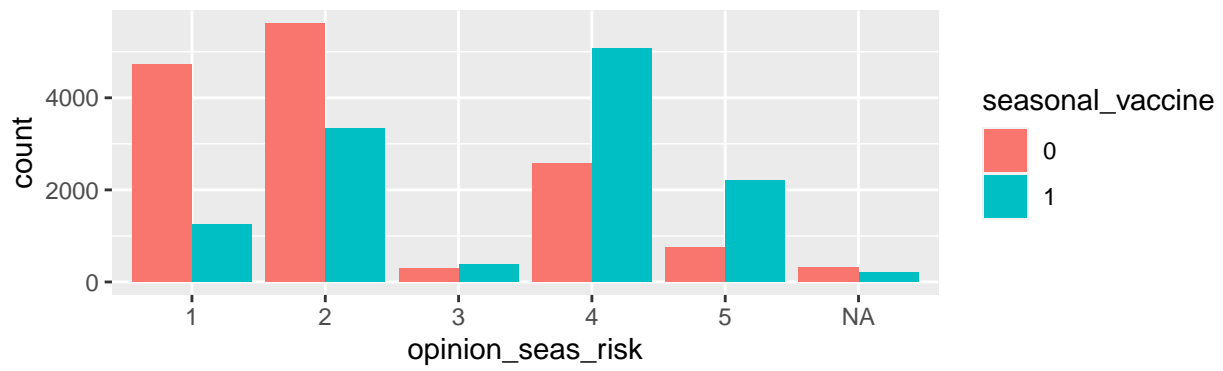
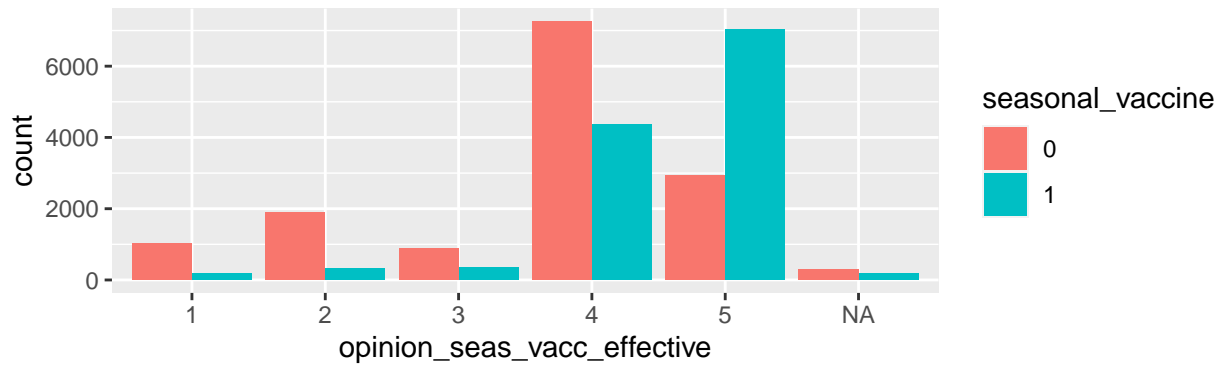


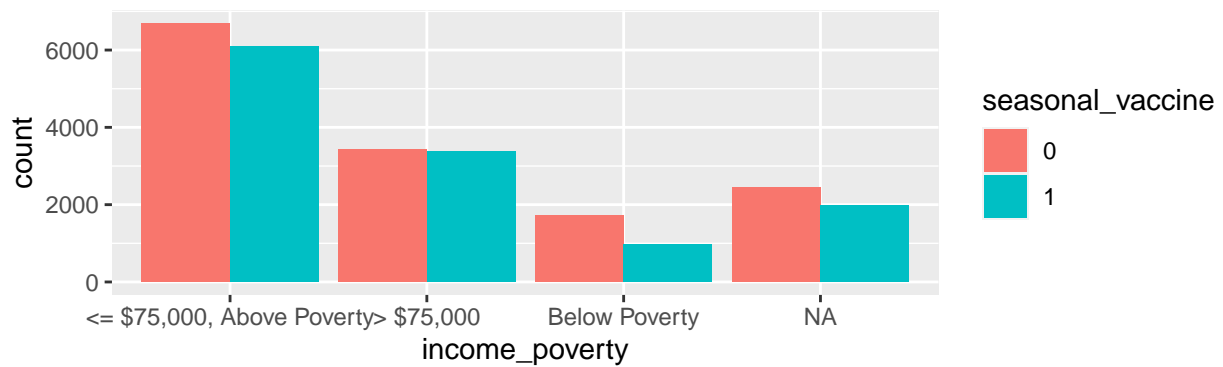
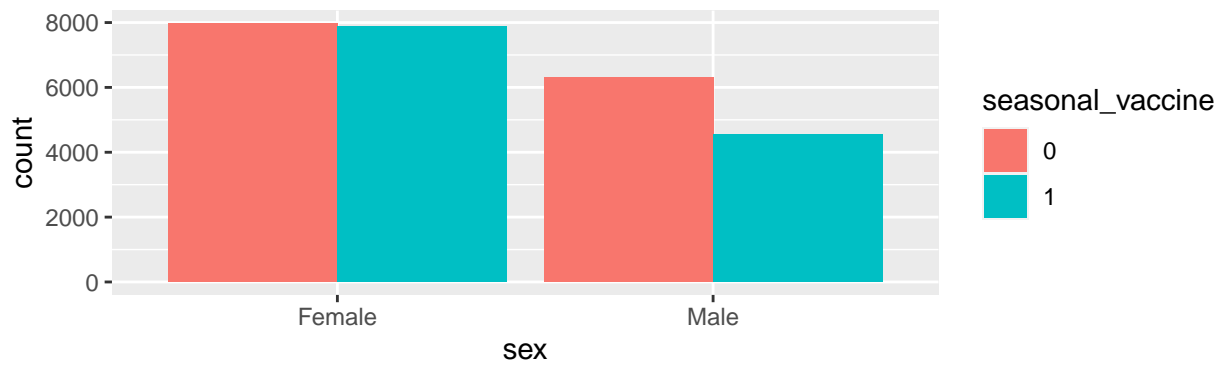
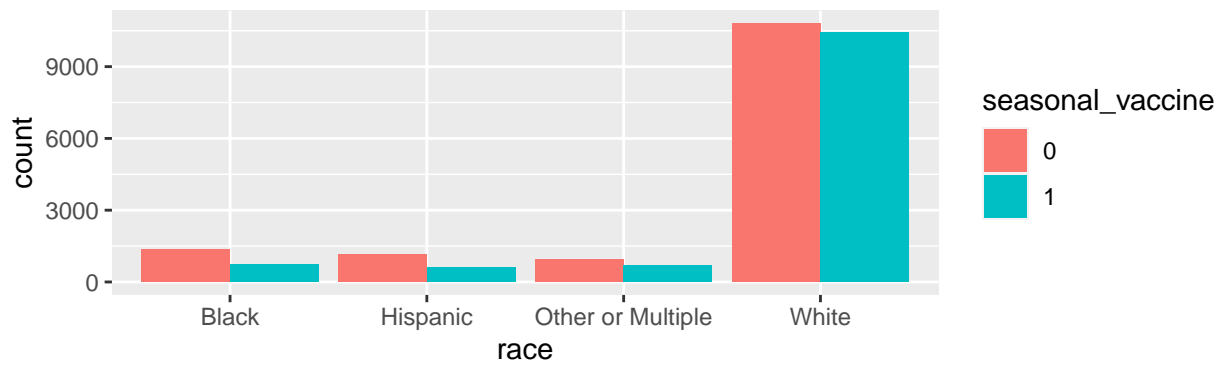
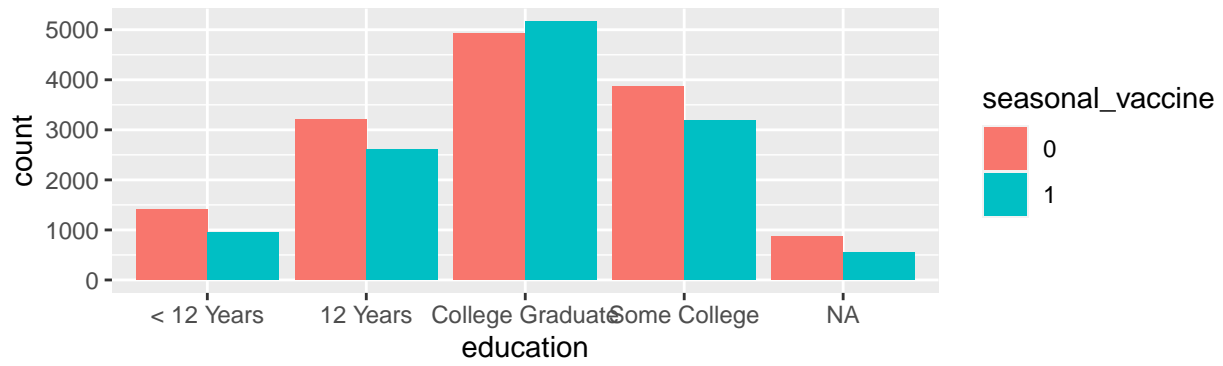


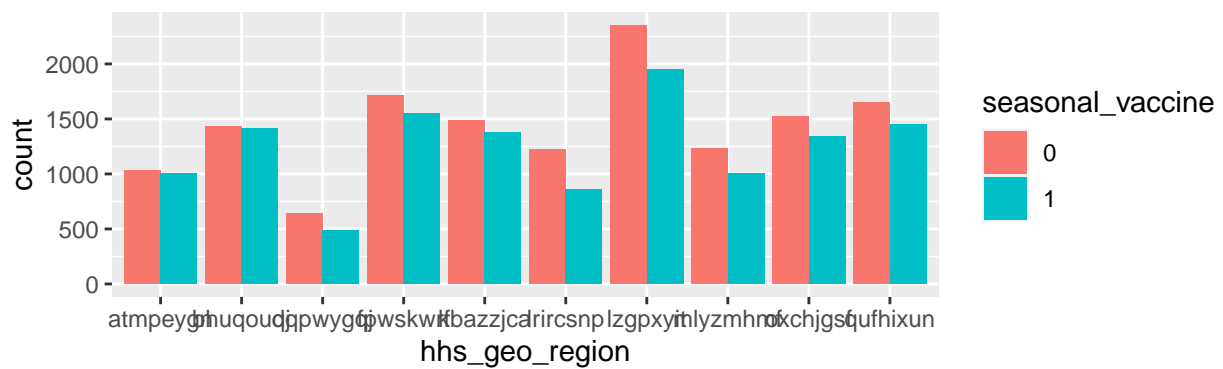
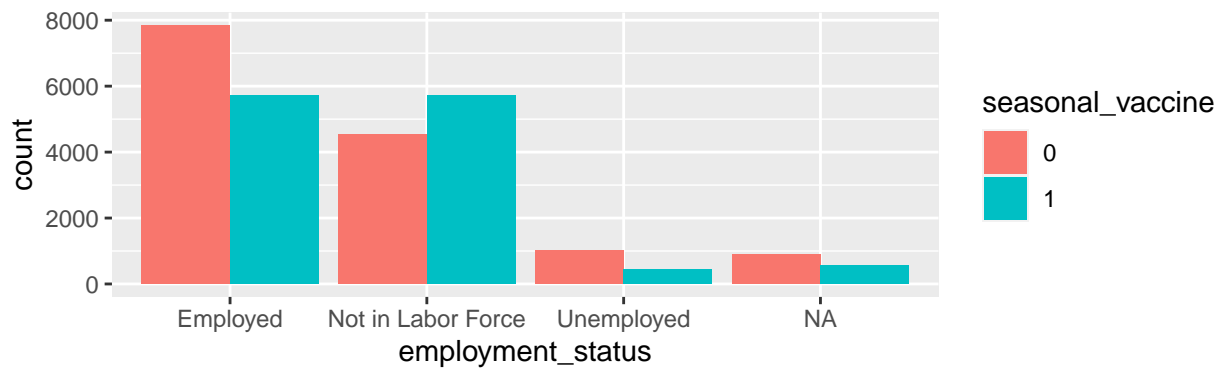
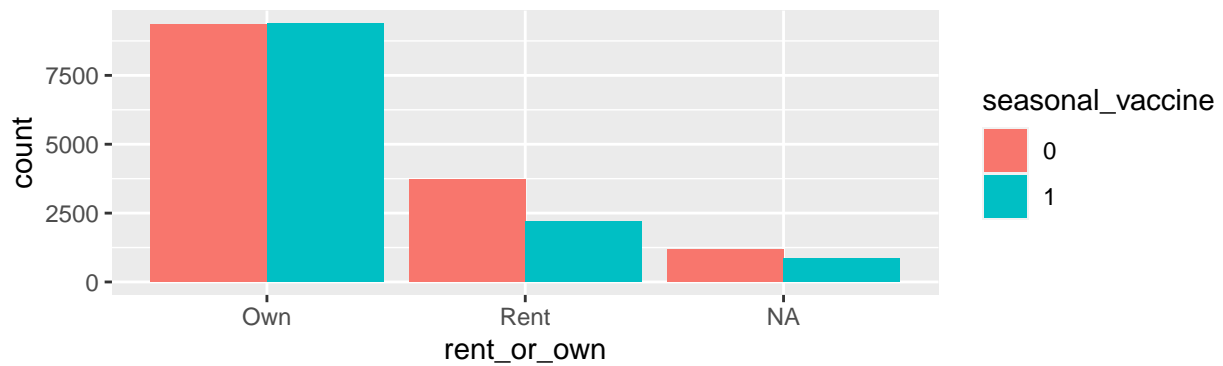
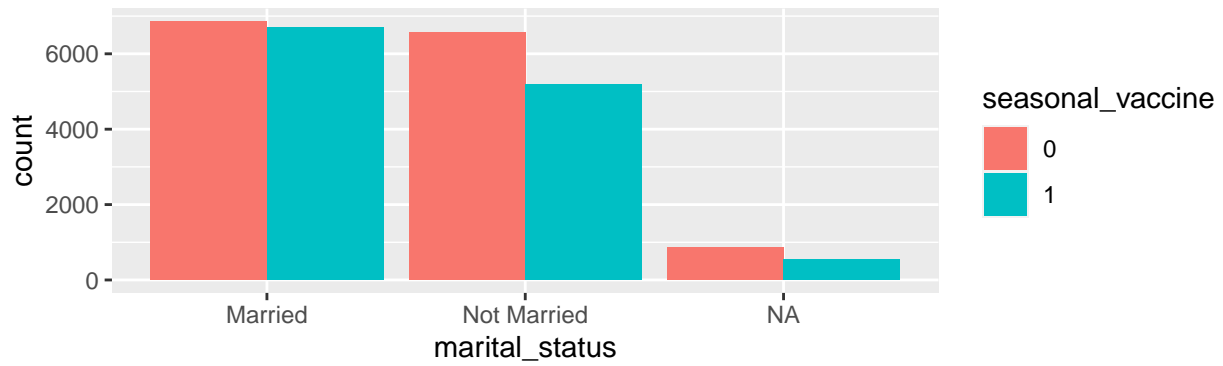


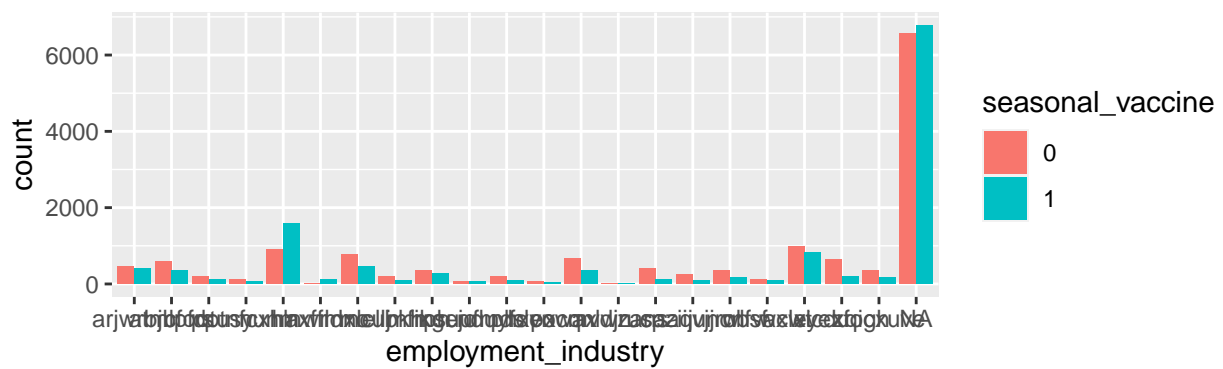
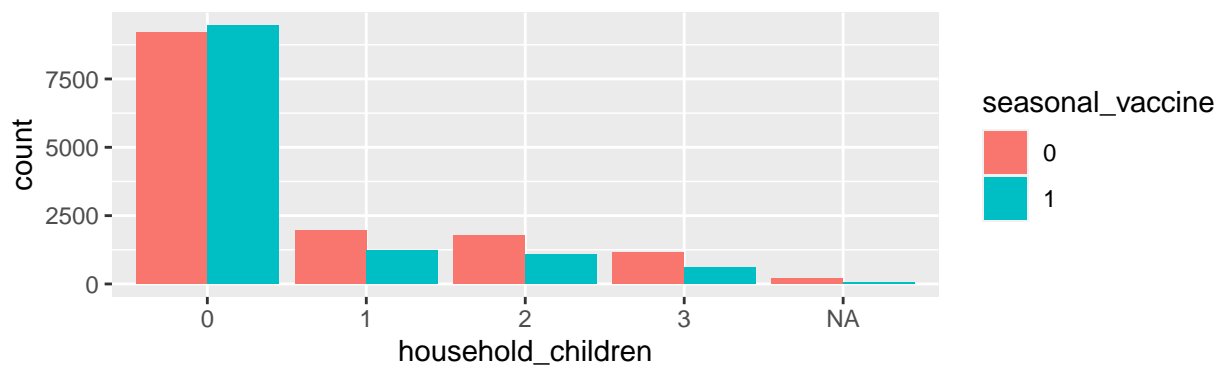
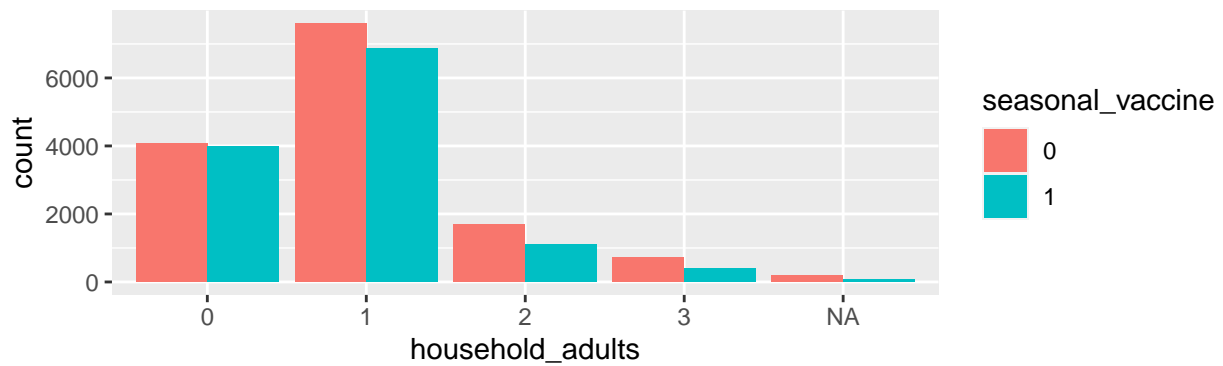






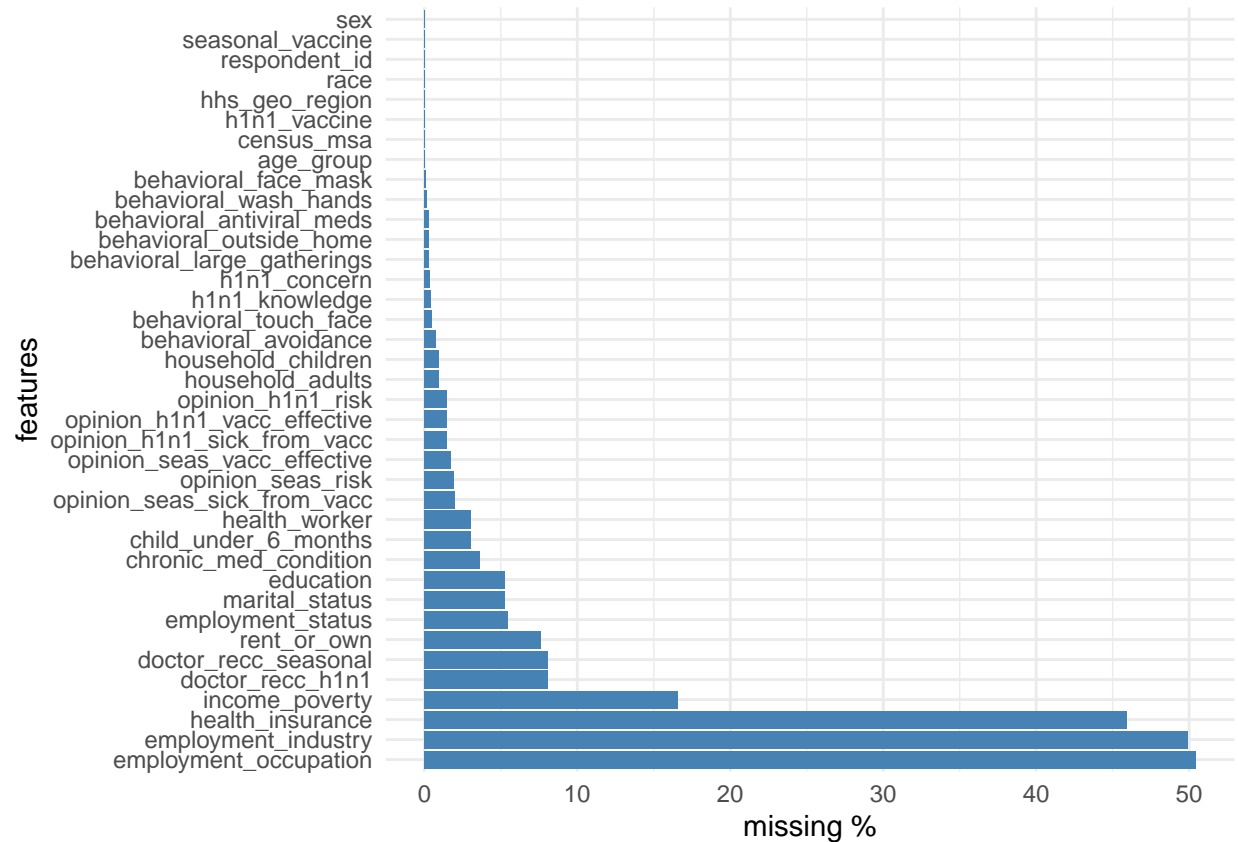






```
df_train %>% summarise_all(funs(sum(is.na(.)/n())*100)) %>%
  gather(key="feature", value="missing_pct") %>%
  ggplot(aes(x=reorder(feature,-missing_pct),y=missing_pct)) +
  geom_bar(stat="identity", fill="steelblue")+
  labs(y = "missing %", x = "features") +
```

```
coord_flip() +  
theme_minimal()
```



```
chisq_values = c()  
uncert_values = c()  
target_values = c()  
  
for(target in targets){  
  for(feature in features){  
    contingency = df_train %>%  
      select(target, feature) %>%  
      table()  
  
    p.value = chisq.test(contingency)$p.value  
  
    target_values = c(target_values, target)  
    chisq_values = c(chisq_values, p.value)  
    uncert_values = c(uncert_values, UncertCoef(contingency))  
  }  
}  
  
results = tibble(feature = rep(features, 2),  
                 target = target_values,  
                 p.value = chisq_values,  
                 uncertainty = uncert_values)
```



```

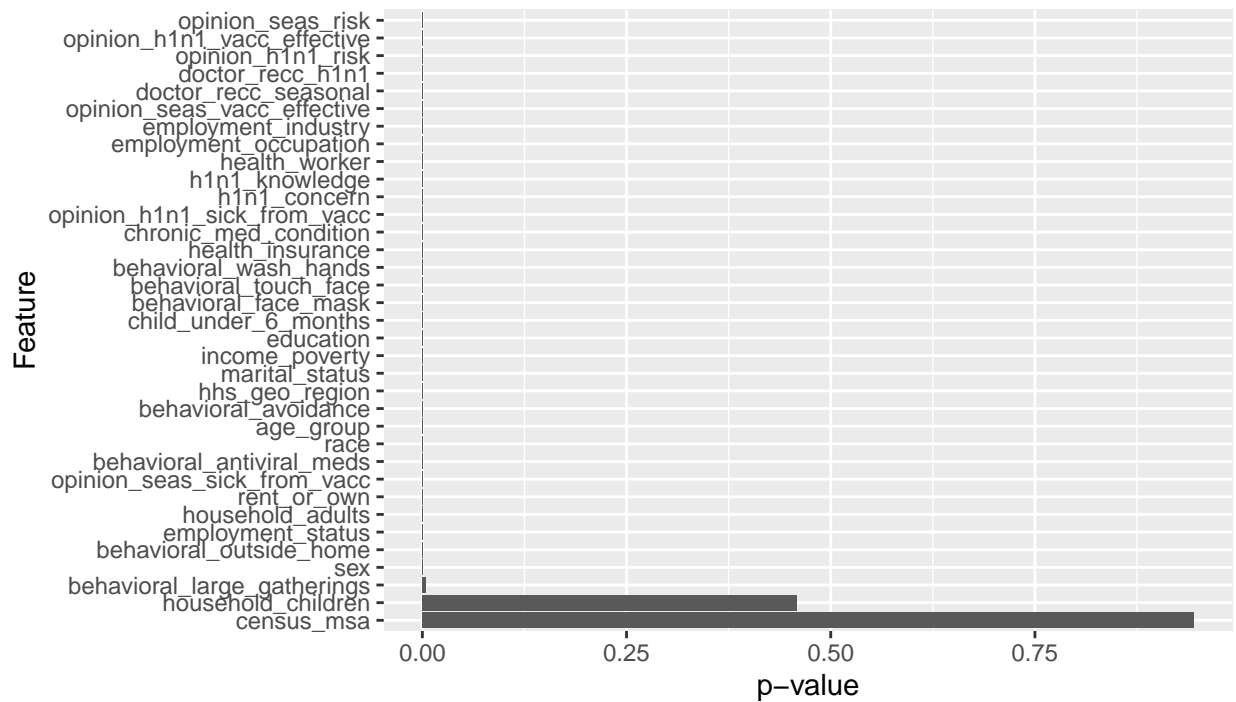
for(target_value in targets){
  p1 = results %>%
    filter(target == target_value) %>%
    filter(p.value < 1) %>%
    ggplot(aes(x=reorder(feature, -p.value), y=p.value)) +
    geom_bar(stat='identity') +
    xlab("Feature") +
    ylab("p-value") +
    ggtitle(sprintf("Chi-squared test p-values for target: %s", target)) +
    coord_flip()

  p2 = results %>%
    filter(target == target_value) %>%
    filter(p.value < 1) %>%
    ggplot(aes(x=reorder(feature, uncertainty), y=uncertainty)) +
    geom_bar(stat='identity') +
    xlab("Feature") +
    ylab("Uncertainty Coefficient") +
    ggtitle(sprintf("Uncertainty Coefficient for target: %s", target)) +
    coord_flip()

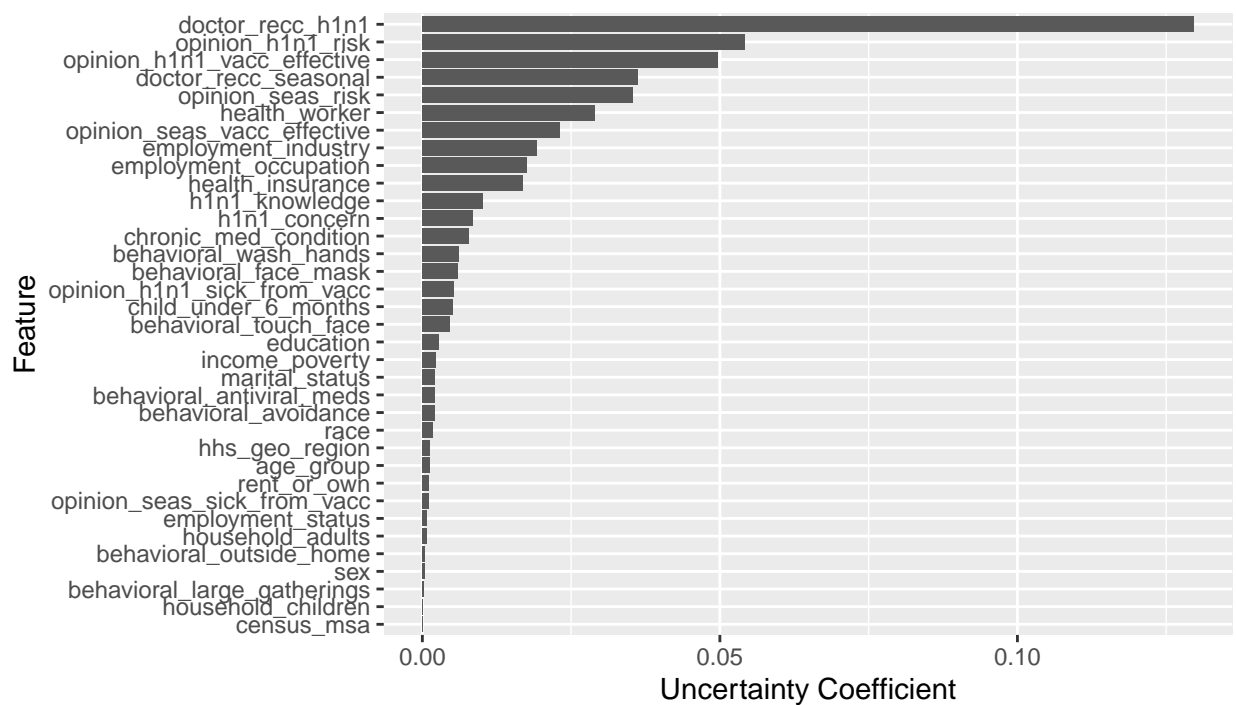
  # print(p1)
  # print(p2)
  do.call("grid.arrange", c(list(p1, p2), nrow=2))
}

```

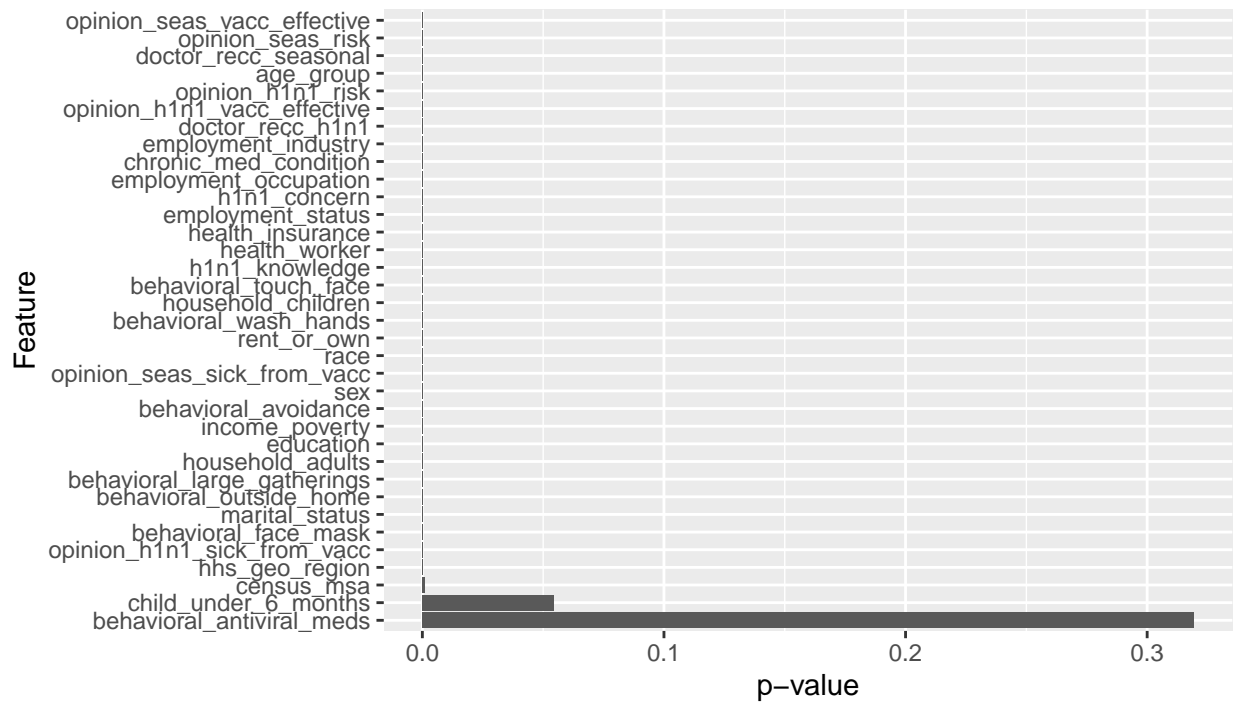
Chi-squared test p-values for target: seasonal_vaccine



Uncertainty Coefficient for target: seasonal_vaccine



Chi-squared test p-values for target: seasonal_vaccine



Uncertainty Coefficient for target: seasonal_vaccine

