# Data Project, Part 3

Emma Collo and Adriel Vijuan

May 06, 2024

**1. [2 marks] The first part of our PPDAC framework is to identify the problem you are addressing with these data. State the question you are trying to answer and let us know what type of question this is in terms of the PPDAC framework. A question statement should be as specific as possible. For example: Do students who regularly get 8 hours of sleep have fewer visits to the health center? This question is an example of an etiologic or causal question.**

Does air pollution affect asthma rates in children? This is an etiologic problem.

**2. [2 marks] Why is this question interesting or important? You could talk here about how existing data/studies suggest this might be important, how the findings might make an impact, how the findings might be used, or why you are personally interested in this question.**

Understanding the correlation between air pollution and asthma rates would provide insight on the health impacts of pollution on vulnerable populations – particularly children. In the context of public health, it could also inform future policies targeted at reducing air pollution levels.Our personal interest in this question is rooted in our experiences of having family members that are asthmatic and are from relatively poor AQI regions.

**3. [2 marks] What is the target population for your project?  Why was this target chosen? (i.e., what was your rationale for wanting to answer this question in this specific population?)**

Our target population is California because the state contains a socioeconomically and ethnically diverse community to help us better understand how confounding variables like income does/does not affect the prevalence of asthma.

**4. [2 marks] What is the sampling frame used to collect the data you are using? It may be helpful here to read any protocol papers, trial registration records, '.Readme' files or documentation that are associated with your dataset. If you have trouble identifying how the records/individuals were sampled, confirm with your supporting GSI that your dataset will be usable for the purposes of the class. Describe why you think this sampling strategy is appropriate for your question. To what group(s) would you feel comfortable generalizing the findings of your study and why?**

Sampling frame: Individuals in California counties, ages 0-17 and 18+ from 2015 to 2020. Who we could generalize the findings to: Individuals 0-17 and 18+ living in regions with similar prevalance of pollutants as California.

**5. [2 marks] Write a brief description (1-4 sentences) of the source and contents of your dataset. Provide a URL to the original data source if applicable. If not (e.g., the data came from your internship), provide 1-2 sentences saying where the data came from. If you completed a web form to access the data and selected a subset, describe these steps (including any options you selected) and the date you accessed the data.**

https://data.chhs.ca.gov/dataset/asthma-prevalence/resource/a440b99b-ccc6-473c-bea1-2baf36b05dbe Accessed: 2/26. This source is from the California Department of Public Health and contains the prevalence of asthma by county, age, 95% confidence interval, year, and grouped county. https://www.lung.org/research/sota/city-rankings/states/california Accessed: 2/26. This source is from the American Lung Association and contains a comprehensive list of counties with their grade and weighted averages of particle pollution https://www.lung.org/research/sota/city-rankings/states/california * Accessed: 2/26. This source is also from the American Lung Association and contains a list of counties and includes risk factors like asthma, poverty, non-White, and pregnancy. https://www.cdph.ca.gov/Programs/CCDPHP/DEODC/EHIB/CPE/Pages/CaliforniaBreathingCountyAsthmaProfiles.aspx Accessed: 2/26. Source: California Department of Public Health. Asthma Prevalence by county based on diagnoses by healthcare providers. Grouped by age.

**6.** **[1 mark] Write code below to import your data into R. Assign your dataset to an object. Make sure to include and annotate this code in your submission (you can use a # to comment out regular text within code chunks to annotate).**

```r
library(readr)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```r
library(ggplot2)
library(readxl)
library(dplyr)
library(readr)
library(broom)
theme_set(theme_minimal())
```

```r
asthma_population<-read_excel("current_asthma.xlsx")

asthma_population <- asthma_population %>%
  rename(
    CURRENT_PREVALENCE = `CURRENT PREVALENCE`,
    AGE_GROUP = `AGE GROUP`
  )
```

**7. [3 marks] Write code in R (included in your submission with annotation) to answer the following questions: i) What are the dimensions of the dataset?**

```
dim(asthma_population)
```

```
## [1] 1239    8
```

ii) What are the variable names of the variables in your dataset?

```
names(asthma_population)
```

```
## [1] "COUNTY"               "YEARS"
## [3] "STRATA"               "AGE_GROUP"
## [5] "CURRENT_PREVALENCE"   "95% CONFIDENCE INTERVAL"
## [7] "COUNTIES GROUPED"     "COMMENT"
```

iii) Print the first six rows of the dataset. 3

```
head(asthma_population)
```

```
## # A tibble: 6 x 8
##    COUNTY     YEARS   STRATA AGE_GROUP CURRENT_PREVALENCE 95% CONFIDENCE INTER~1
##    <chr>      <chr>   <chr>  <chr>                  <dbl> <chr>
## 1 California 2015-20~ Total~ All ages              0.0870 (8.1-9.3)
## 2 Alameda    2015-20~ Total~ All ages              0.0905 (5.9-12.2)
## 3 Alpine     2015-20~ Total~ All ages              0.093  (4.1-14.6)
## 4 Amador     2015-20~ Total~ All ages              0.093  (4.1-14.6)
## 5 Butte      2015-20~ Total~ All ages              0.0943 (3.8-15.1)
## 6 Calaveras  2015-20~ Total~ All ages              0.093  (4.1-14.6)
## # i abbreviated name: 1: `95% CONFIDENCE INTERVAL`
## # i 2 more variables: `COUNTIES GROUPED` <chr>, COMMENT <chr>
```

**8. [2 marks]** Use the data to demonstrate a data visualization skill we have covered during Part I of the course. Choose a visualization relevant to your stated problem. Include your code in your submission. For example, you could visualize the distribution of our outcome with a histogram, or use a bar graph to represent the distribution of your exposure variable.

```
CP_data<-asthma_population%>%mutate(CP_DATA_100=100*CURRENT_PREVALENCE)
CP_data
```

```
## # A tibble: 1,239 x 9
##     COUNTY        YEARS STRATA AGE_GROUP CURRENT_PREVALENCE 95% CONFIDENCE INTER~1
##     <chr>         <chr> <chr>  <chr>                  <dbl> <chr>
##  1 California    2015~ Total~ All ages              0.0870 (8.1-9.3)
##  2 Alameda       2015~ Total~ All ages              0.0905 (5.9-12.2)
##  3 Alpine        2015~ Total~ All ages              0.093  (4.1-14.6)
##  4 Amador        2015~ Total~ All ages              0.093  (4.1-14.6)
##  5 Butte         2015~ Total~ All ages              0.0943 (3.8-15.1)
##  6 Calaveras     2015~ Total~ All ages              0.093  (4.1-14.6)
##  7 Colusa        2015~ Total~ All ages              0.0733 (2.7-12.0)
##  8 Contra Costa  2015~ Total~ All ages              0.121  (5.2-19.1)
##  9 Del Norte     2015~ Total~ All ages              0.0722 (1.2-13.3)
## 10 El Dorado     2015~ Total~ All ages              0.121  (2.9-21.2)
## # i 1,229 more rows
## # i abbreviated name: 1: '95% CONFIDENCE INTERVAL'
## # i 3 more variables: 'COUNTIES GROUPED' <chr>, COMMENT <chr>,
## #   CP_DATA_100 <dbl>
```

```
CP_data_log<-CP_data%>%mutate(log_CP=log(CP_DATA_100))
CP_data_log
```

```
## # A tibble: 1,239 x 10
##     COUNTY        YEARS STRATA AGE_GROUP CURRENT_PREVALENCE 95% CONFIDENCE INTER~1
##     <chr>         <chr> <chr>  <chr>                  <dbl> <chr>
##  1 California    2015~ Total~ All ages              0.0870 (8.1-9.3)
##  2 Alameda       2015~ Total~ All ages              0.0905 (5.9-12.2)
##  3 Alpine        2015~ Total~ All ages              0.093  (4.1-14.6)
##  4 Amador        2015~ Total~ All ages              0.093  (4.1-14.6)
##  5 Butte         2015~ Total~ All ages              0.0943 (3.8-15.1)
##  6 Calaveras     2015~ Total~ All ages              0.093  (4.1-14.6)
##  7 Colusa        2015~ Total~ All ages              0.0733 (2.7-12.0)
##  8 Contra Costa  2015~ Total~ All ages              0.121  (5.2-19.1)
##  9 Del Norte     2015~ Total~ All ages              0.0722 (1.2-13.3)
## 10 El Dorado     2015~ Total~ All ages              0.121  (2.9-21.2)
## # i 1,229 more rows
## # i abbreviated name: 1: '95% CONFIDENCE INTERVAL'
## # i 4 more variables: 'COUNTIES GROUPED' <chr>, COMMENT <chr>,
## #   CP_DATA_100 <dbl>, log_CP <dbl>
```

```
asthma_17<-CP_data_log%>%filter(AGE_GROUP!="All ages",AGE_GROUP!="18-64 years",AGE_GROUP!="18+ years",AG
asthma_17
```

```
## # A tibble: 354 x 10
##     COUNTY        YEARS STRATA AGE_GROUP CURRENT_PREVALENCE 95% CONFIDENCE INTER~1
```

```
##    <chr>          <chr> <chr>  <chr>                      <dbl> <chr>
##  1 California     2015~ Age g~ 0-4 years               0.0445 (2.5-6.4)
##  2 Alameda        2015~ Age g~ 0-4 years                   NA <NA>
##  3 Alpine         2015~ Age g~ 0-4 years                   NA <NA>
##  4 Amador         2015~ Age g~ 0-4 years                   NA <NA>
##  5 Butte          2015~ Age g~ 0-4 years                   NA <NA>
##  6 Calaveras      2015~ Age g~ 0-4 years                   NA <NA>
##  7 Colusa         2015~ Age g~ 0-4 years                   NA <NA>
##  8 Contra Costa   2015~ Age g~ 0-4 years                   NA <NA>
##  9 Del Norte      2015~ Age g~ 0-4 years                   NA <NA>
## 10 El Dorado      2015~ Age g~ 0-4 years                   NA <NA>
## # i 344 more rows
## # i abbreviated name: 1: `95% CONFIDENCE INTERVAL`
## # i 4 more variables: `COUNTIES GROUPED` <chr>, COMMENT <chr>,
## #   CP_DATA_100 <dbl>, log_CP <dbl>
```

```r
asthma_17_county<-asthma_17%>%filter(CURRENT_PREVALENCE>0,COUNTY!="Sutter")
```

```r
ggplot(asthma_17_county,aes(x=COUNTY,y=log_CP))+geom_bar(stat="identity",aes(col=AGE_GROUP),position="d
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>
```

```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>
```

### Current Prevalence of Asthma in individuals under 17yrs 2015–2020.



```
asthma_youth_v_old<-CP_data_log%>%filter(AGE_GROUP!="All ages",AGE_GROUP!="65+ years",AGE_GROUP!="0-4 ye
asthma_youth_v_old
```

```
## # A tibble: 225 x 10
##    COUNTY      YEARS STRATA AGE_GROUP CURRENT_PREVALENCE 95% CONFIDENCE INTER~1
##    <chr>       <chr> <chr>  <chr>                  <dbl> <chr>
## 1 California  2015~ Child~ 0-17 yea~              0.101 (8.3-11.9)
## 2 Alameda     2015~ Child~ 0-17 yea~              0.151 (3.9-26.4)
## 3 Fresno      2015~ Child~ 0-17 yea~              0.107 (0.5-20.8)
## 4 Imperial    2015~ Child~ 0-17 yea~              0.191 (5.6-32.5)
## 5 Kern        2015~ Child~ 0-17 yea~              0.260 (10.8-41.1)
## 6 Los Angeles 2015~ Child~ 0-17 yea~              0.101 (6.0-14.2)
## 7 Orange      2015~ Child~ 0-17 yea~              0.0799 (2.4-13.6)
```

```
##  8 Riverside    2015~ Child~ 0-17 yea~              0.148  (3.1-26.4)
##  9 San Bernard~ 2015~ Child~ 0-17 yea~              0.118  (2.1-21.6)
## 10 San Diego    2015~ Child~ 0-17 yea~              0.0658 (2.4-10.7)
## # i 215 more rows
## # i abbreviated name: 1: '95% CONFIDENCE INTERVAL'
## # i 4 more variables: 'COUNTIES GROUPED' <chr>, COMMENT <chr>,
## #   CP_DATA_100 <dbl>, log_CP <dbl>
```

```r
ggplot(asthma_youth_v_old, aes(x = COUNTY, y = log_CP)) +
  geom_point(aes(color = AGE_GROUP)) +
  labs(
    x = "Counties",
    y = "Current Prevalence (LOG)",
    title = "Current Prevalence of Asthma for Individuals All Ages 2015-2020."
  ) +
  facet_wrap(~AGE_GROUP, nrow = 2) +
  theme(
    axis.text.x = element_text(angle = 90, hjust = 1),
    plot.title = element_text(size = 12)
  )
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>
```

# Current Prevalence of Asthma for Individuals All Ages 2015−2020.

**9. [2 marks] Describe the skill that you are demonstrating and interpret your findings. For example, if you have created a histogram, describe the central tendency, shape of the distribution, etc.**

When looking at the youth asthma cases, Sutter County was an outlier in the dataset because it "statistically unstable." Within youth cases, there is a higher asthma incidence in ages 5-17 compared to 0-4. Additionally when comparing asthma cases between ages, there is a higher distribution of asthma cases in individuals that are 18+.

10. **[1 mark] Include your work for Part I.**

**11. [3 marks] Describe the type of theoretical distribution that is relevant for your data. What type of variable(s) are you investigating (continuous, categorical, ordinal, etc)?**

Current Prevalence – continuous Counties – nomial

**What theoretical distribution that we have talked about would potentially be appropriate to use with these data (Normal, Binomial, Poisson. . . )**

Normal

**Why is this an appropriate model for the data you are studying? (HINT what are the assumptions of this distribution)**

To determine if the data being used follows a normal distribution, we plotted the data on a QQ plot to better visualize normality. The data shows that 0-17yrs and 18+ yrs has a normal distribution because the points are close to the line.

**12.[2 marks] - What are the parameters that define this distribution? - Calculate these parameters for your data.**

```r
asthma17<-asthma_17%>%filter(CURRENT_PREVALENCE>0)
```

```r
asthma_17_summary<-asthma17%>%summarize(mean_asthma17=mean(CURRENT_PREVALENCE),sd_asthma17=sd(CURRENT_PI
```

x= 0.1114921, sd = 0.06089466

**13. [2 marks] Use your outcome data to calculate a probability. Provide an equation (use fpr,a; probability notation) that describes this probability. Note whether this probability is a conditional or a marginal probability. For example if my outcome variable is height in inches, I might calculate the probability that an individual in the dataset has a height of greater than 5 feet 10 inches. P(height>=70inches) = ? This would be a marginal probability.**

P(diagnosed with asthma, < 17yrs.from 2019-2020) marginal

```
asthma_17_2019<-asthma17%>%filter(YEARS!="2017-2018",YEARS!="2015-2016",CURRENT_PREVALENCE>0)
```

```
asthma17_2019<-asthma_17_2019%>%summarize(mean_asthma17_2019=mean(CURRENT_PREVALENCE))
(0.1042367)*(0.1114921)
```

```
## [1] 0.01162157
```

P(diagnosed with asthma, < 17yrs.from 2019-2020) = 0.01162157

P(diagnosed with asthma, 18+ in 2019-2020) marginal

```
asthma18<-CP_data_log%>%filter(AGE_GROUP!="All ages",AGE_GROUP!="65+ years",AGE_GROUP!="0-4 years",AGE_
```

```
asthma18_2019<-asthma18%>%summarize(mean_asthma18=mean(CURRENT_PREVALENCE))
(0.1070552)*(0.1114921)
```

```
## [1] 0.01193581
```

P(diagnosed with asthma, 18+ in 2019-2020) = 0.01193581

**14. [2 marks] What type of variable is your primary exposure of interest? If this variable is a demographic variable (age, gender identity, race/ethnic identity) explain how the categories of this variable are defined and what the rationale is for this (for example if gender identity is being used, is the idea to capture something about biology ie using gender identity as a marker for genetic or phenotypic sex, or as a marker of social exposures). If your data is not from a randomized trial where your exposure of interest was randomly assigned, are there important factors that may have affected how this exposure was distributed?**

Our primary exposure of interest is age. We specifically want to observe if asthma rates in children are increasing over time, and from there we can further investigate the specific county that has high asthma rates. Additionally, we will observe if there is a correlation between high asthma rates in children and adults in the same county.

**15.** **[4 marks] Use your data to create a visualization of your data that begins to explore your research question. Include code in R, a visual of some kind and text interpretation. For example, if you outcome is height of children at age 10 and your predictor variable is exposure to food insecurity in the first year of life, you could provide a histogram of height among children exposed to food insecurity and a separate histogram of height among children not exposed to food insecurity. Make sure you describe your interpretation of the results.**

```r
ggplot(asthma_youth_v_old, aes(x = COUNTY, y = log_CP)) +
  geom_point(aes(color = YEARS)) +
  geom_line(aes(color = YEARS, group = YEARS)) +
  labs(
    x = "Counties",
    y = "Current Prevalence (LOG)",
    title = "Current Prevalence of Asthma for Individuals All Ages Separated by Year."
  ) +
  facet_wrap(~AGE_GROUP, nrow = 2) +
  theme_minimal() +
  theme(
    axis.text.x = element_text(angle = 90, hjust = 1),
    plot.title = element_text(size = 12)
  )
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-17 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>
```

# Current Prevalence of Asthma for Individuals All Ages Separated by Year.



```
desc_youth_v_old<-asthma_youth_v_old%>%arrange(desc(log_CP))

youth<-desc_youth_v_old%>%filter(AGE_GROUP!="18+ years",CURRENT_PREVALENCE>0)

ggplot(youth, aes(x = COUNTY, y = log_CP)) +
  geom_point(aes(color = YEARS)) +
  geom_line(aes(color = YEARS, group = YEARS)) +
  geom_smooth(method = "loess", se = FALSE) +
  labs(
    x = "Counties",
    y = "Current Prevalence (LOG)",
    title = "Current Prevalence of Asthma 0-17yrs. in 2015"
  ) +
  facet_wrap(~YEARS, nrow = 3) +
  theme_minimal() +  #
  theme(
    axis.text.x = element_text(angle = 90, hjust = 1),
    plot.title = element_text(size = 12)
  )
```

```
## `geom_smooth()` using formula = 'y ~ x'

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>
```
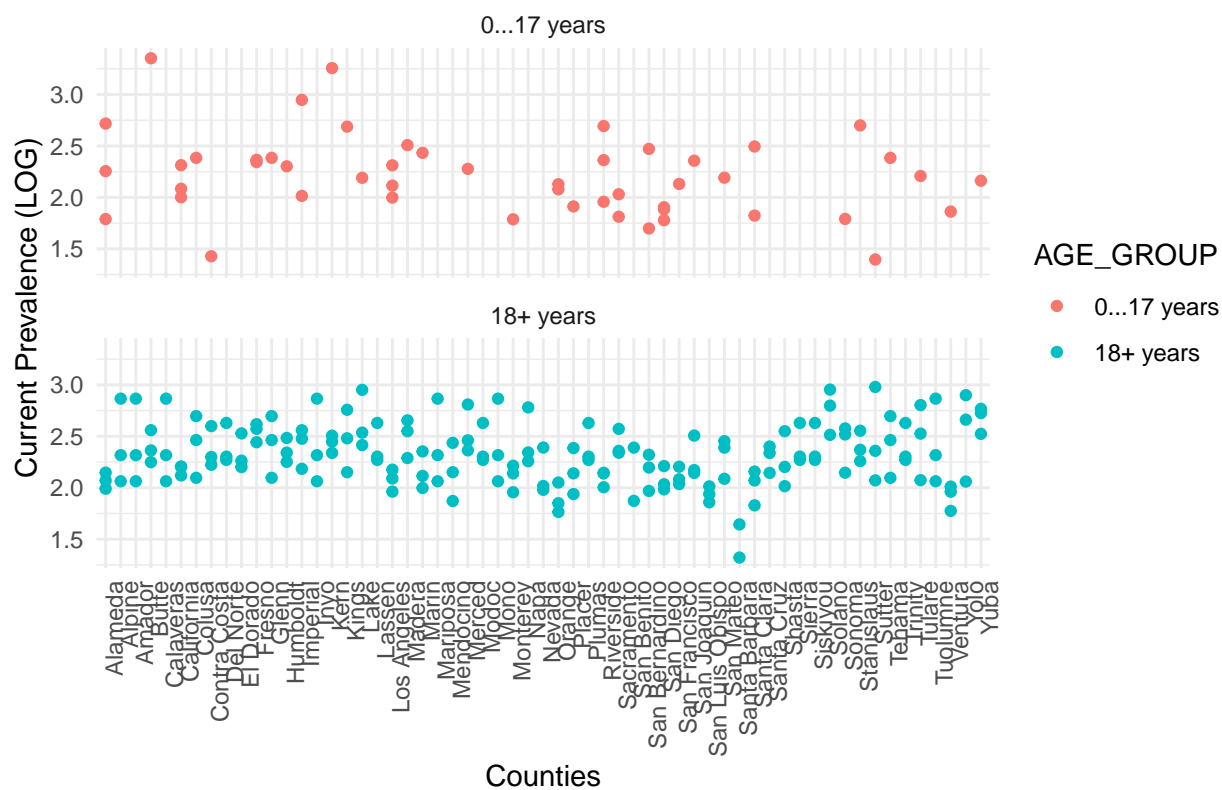
```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>
```

```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>
```

## Current Prevalence of Asthma 0–17yrs. in 2015



```
old<-asthma_youth_v_old%>%filter(AGE_GROUP!="0-17 years",CURRENT_PREVALENCE>0)
```

```
ggplot(old, aes(x = COUNTY, y = log_CP)) +
  geom_point(aes(color = YEARS)) +
  geom_line(aes(color = YEARS, group = YEARS)) +
  labs(
    x = "Counties",
    y = "Current Prevalence (LOG)",
    title = "Current Prevalence of Asthma 18+yrs. in 2015"
  ) +
  facet_wrap(~YEARS, nrow = 3) +
  theme_minimal() +
  theme(
    axis.text.x = element_text(angle = 90, hjust = 1),
    plot.title = element_text(size = 12)
  )
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>
```
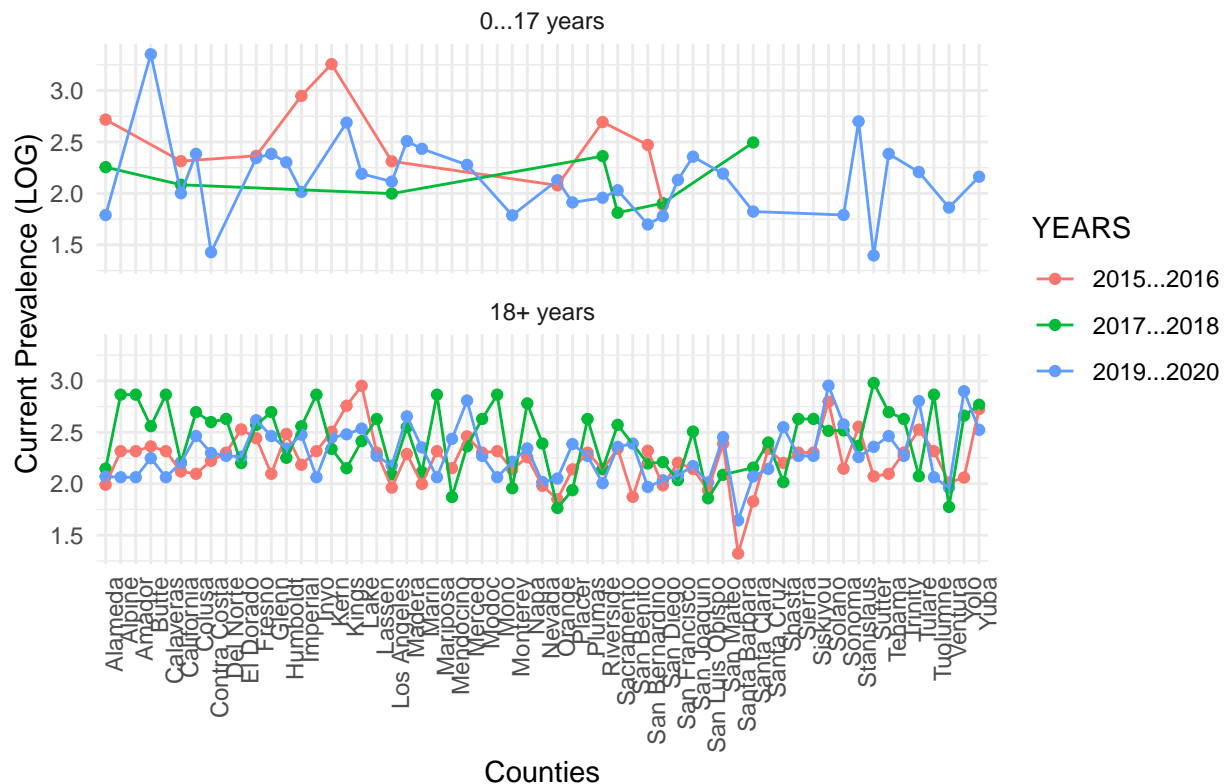
```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>
```
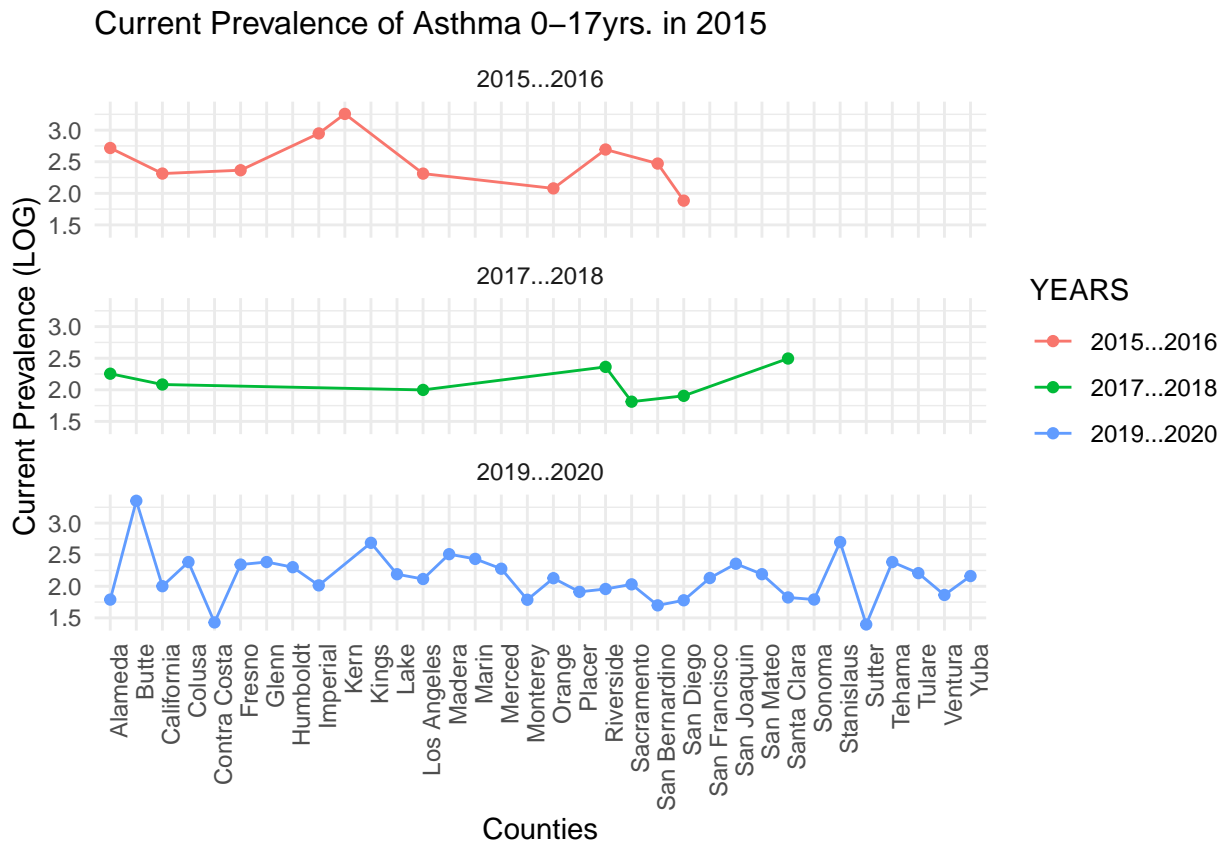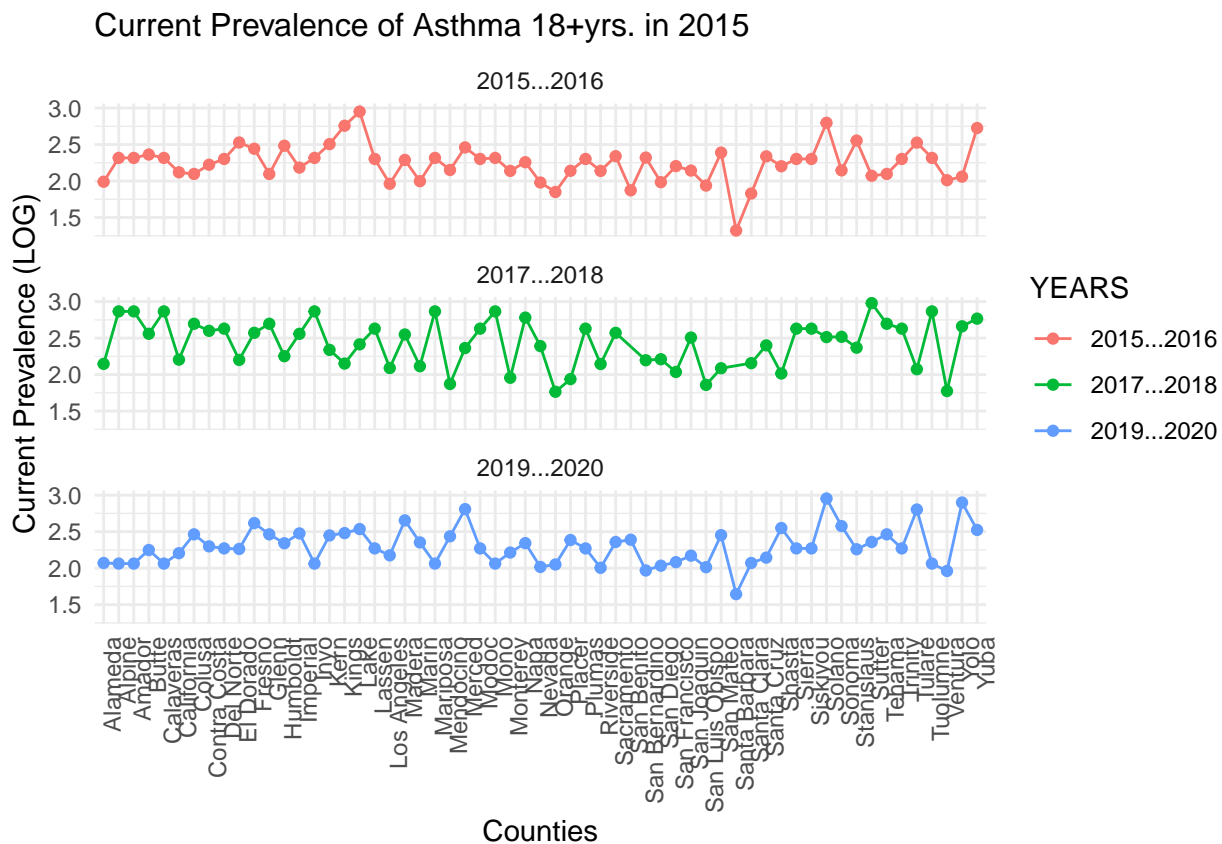
```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>
```

## Current Prevalence of Asthma 18+yrs. in 2015



```
ggplot(asthma_17_county, aes(x = COUNTY, y = log_CP)) +
  geom_point(aes(color = AGE_GROUP)) +
  labs(
    x = "Counties",
    y = "Current Prevalence (LOG)",
    title = "Current Prevalence of Asthma for 0-17 yrs. from 2015-2020"
  ) +
  facet_wrap(~YEARS, nrow = 3) +
  theme_minimal() +
  theme(
    axis.text.x = element_text(angle = 90, hjust = 1),
    plot.title = element_text(size = 12),
    axis.text.y = element_text(size = 12)
  )
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2019-2020' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <e2>
```
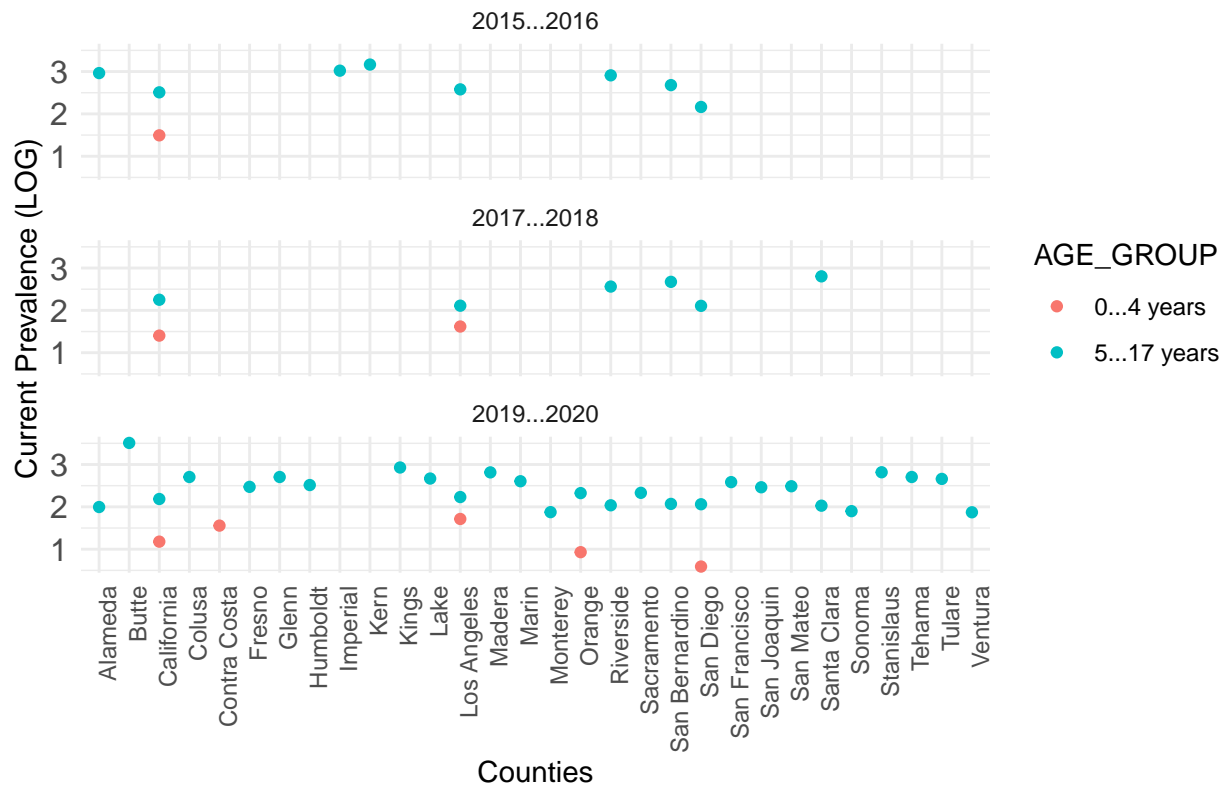
```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2017-2018' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '2015-2016' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>
```
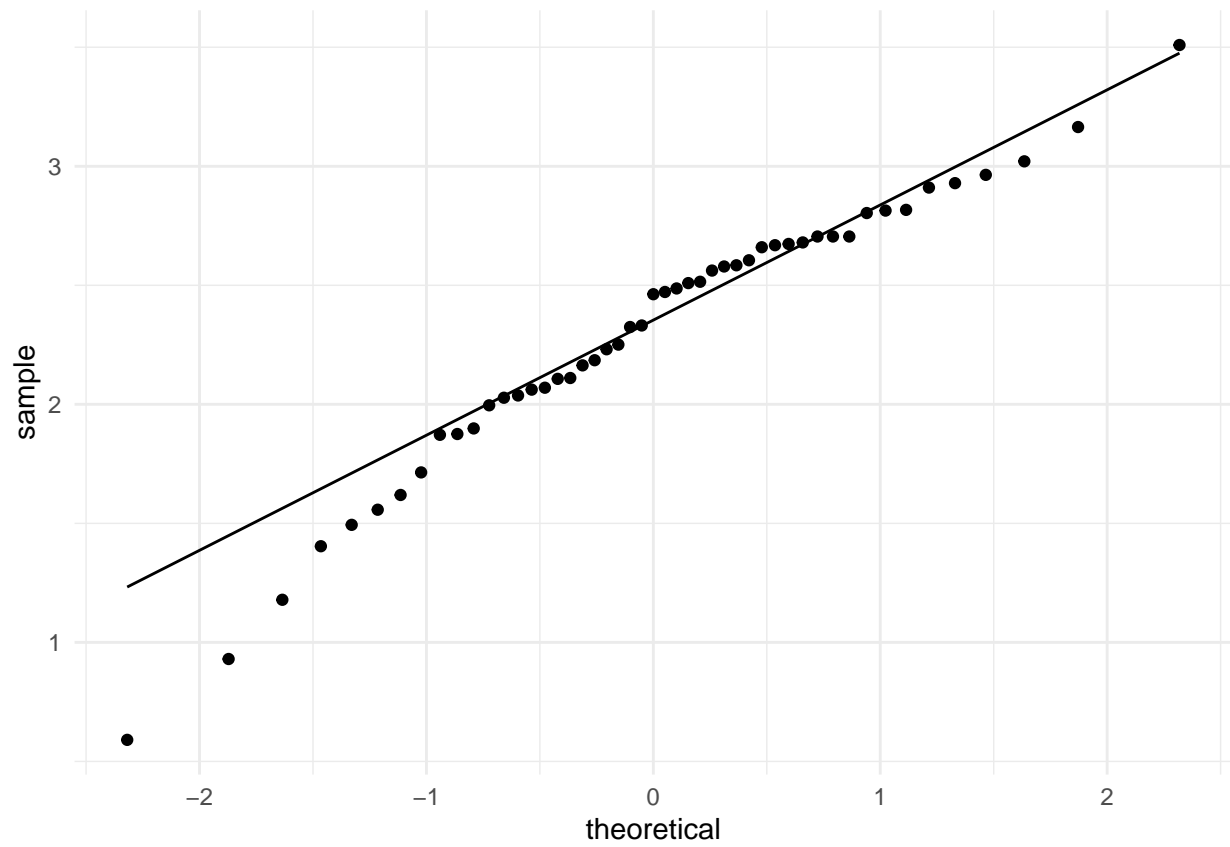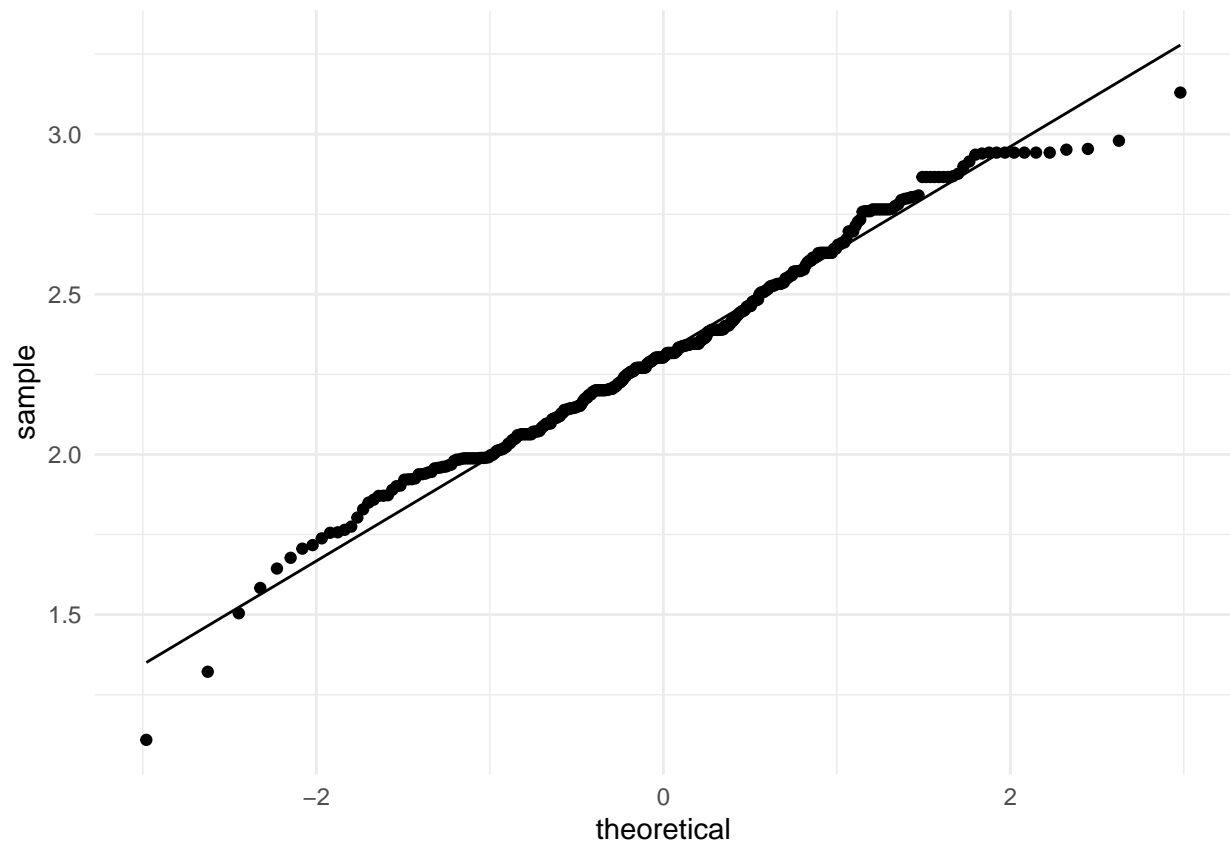
# Current Prevalence of Asthma for 0–17 yrs. from 2015−2020



```
ggplot(asthma_17_county,aes(sample=log_CP))+stat_qq_line()+stat_qq()
```

```
ggplot(asthma18,aes(sample=log_CP))+stat_qq_line()+stat_qq()
```

**16. [1 mark] Include parts I and II of your project.**

See above.

**17. [2 marks] Identify a statistical test to apply to your data (must be a statistical test that we cover in part III of the course). Name the statistical test you have chosen and explain why this is the appropriate test for these data. (for example, if I have a pre and post intervention measure of morning sleepyness that is quantitative, I might choose a paired t test, because the paired t test is appropriate for continuous outcome data in 2 groups that are inherently related)**

We chose to use a regression inference to explore the effects of AQI levels and age groups on asthma rates and this method allows for the examination of continuous response variables (asthma rates) and how they vary with both categorical predictors (AQI levels, age groups).

```
asthma_pt3<-CP_data_log%>%filter(YEARS!="2015-2016",YEARS!="2017-2018",AGE_GROUP!="All ages", AGE_GROUP
```

**18. [2 marks] What assumptions are required by the testing method you chose? Are these assumptions met by your data? How did you assess this? For example, one of the assumptions of the t-test is that the data are normally distributed, so you might choose to assess this with a histogram, or a q-q plot.**

Some assumptions we concluded were: Asthma rates in different counties are assumed to be independent of each other Relationship between predictors (AQI levels/Current Prevalence, age) and asthma rates is linear Residuals are normally distributed. The diagnostic plots used in regression analysis help to evaluate the fit and assumptions of the model. Plot a (Residuals vs Fitted) and Plot c (Fitted Values vs Residuals) check for non-linearity, inconsistent variance, or outliers by showing residuals against predicted values, while Plot b (Normal Q-Q Plot of Residuals) assesses the normality of residuals by comparing them to a theoretical normal distribution. Plot d (Boxplot of Residuals) summarizes the distribution of residuals, indicating model fit by showing their spread and identifying potential outliers.

```r
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----------------------- tidyverse 2.0.0 --
## v forcats   1.0.0     v stringr   1.5.1
## v lubridate 1.9.3     v tibble    3.2.1
## v purrr     1.0.2     v tidyr     1.3.1
## -- 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 error
```

```r
asthma_data <- read_excel("current_asthma.xlsx") %>%
  rename(
    AGE_GROUP = `AGE GROUP`,
    CURRENT_PREVALENCE = `CURRENT PREVALENCE`
  ) %>%
  filter(YEARS != "2015-2016", YEARS != "2017-2018", CURRENT_PREVALENCE > 0)
```

```r
asthma_lm <- lm(log(CURRENT_PREVALENCE) ~ AGE_GROUP, data = asthma_data)
```

```r
augmented_data <- augment(asthma_lm)

# calculate standardized residuals if not present
augmented_data$.stdresid <- rstandard(asthma_lm)

# plot a: scatter plot with regression line and residuals
plot_a <- ggplot(augmented_data, aes(x = .fitted, y = .resid)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE, color = "red", linetype = "dashed") +
  labs(title = "Plot a: Residuals vs Fitted with Regression Line",
       x = "Fitted Values", y = "Residuals") +
  geom_hline(yintercept = 0, linetype = "dotted")

# plot b: qq plot of residuals
plot_b <- ggplot(augmented_data, aes(sample = .stdresid)) +
  stat_qq() +
  stat_qq_line() +
  labs(title = "Plot b: Normal Q-Q Plot of Residuals",
       x = "Theoretical Quantiles", y = "Standardized Residuals")
```

```r
# plot c: fitted values vs residuals
plot_c <- ggplot(augmented_data, aes(x = .fitted, y = .resid)) +
  geom_point() +  # this ensures it's a scatter plot
  geom_hline(yintercept = 0, linetype = "dotted") +
  labs(title = "Plot c: Fitted Values vs Residuals",
       x = "Fitted Values", y = "Residuals")


# plot d: boxplot of the distribution of y and residuals
plot_d <- ggplot(augmented_data, aes(x = factor(1), y = .resid)) +
  geom_boxplot() +
  labs(title = "Plot d: Boxplot of Residuals",
       x = "", y = "Residuals")

# display plots and sse
library(gridExtra)
```

```
##
## Attaching package: 'gridExtra'

## The following object is masked from 'package:dplyr':
##
##     combine
```

```r
grid.arrange(plot_a, plot_b, plot_c, plot_d, nrow = 2)
```

```
## `geom_smooth()` using formula = 'y ~ x'
```

Plot a: Residuals vs Fitted with Regression Line  Plot b: Normal Q–Q Plot of Residuals

Plot c: Fitted Values vs Residuals  Plot d: Boxplot of Residuals

19. **[2 marks] Clearly state the null and alternative hypotheses for your test.**

- **Null Hypothesis (H0):** There is no difference in asthma rates in age groups (0-4, 5-17, 18-64, and 65+ years) in 2019 to 2020. -**Alternative Hypothesis (Ha):** At least one age group has a relationship with the prevalence of asthma rates in 2019 to 2020.

**20.** **[2 marks]** Conduct the statistical test. Include the R code you used to generate your results. Annotate your code to help us follow your reasoning.

```r
# Filter the data to exclude certain years and age groups, ensuring the prevalence is above zero
asthma_pt3 <- CP_data_log %>%
  filter(
    YEARS != "2015-2016", YEARS != "2017-2018",
    AGE_GROUP != "All ages", AGE_GROUP != "18+ years", AGE_GROUP != "0-17 years",
    CURRENT_PREVALENCE > 0
  )

# Conduct a linear regression analysis where the log of current prevalence is explained by age groups
asthma_lm <- lm(log_CP ~ AGE_GROUP, data = asthma_pt3)

# Display the first few rows to inspect the filtered dataset
head(asthma_pt3)
```

```
## # A tibble: 6 x 10
##   COUNTY      YEARS  STRATA AGE_GROUP CURRENT_PREVALENCE 95% CONFIDENCE INTER~1
##   <chr>       <chr>  <chr>  <chr>                  <dbl> <chr>
## 1 California  2019-~ Age g~ 0-4 years            0.0325  (2.2-4.3)
## 2 Contra Costa 2019-~ Age g~ 0-4 years           0.0475  (0.3-9.2)
## 3 Los Angeles 2019-~ Age g~ 0-4 years            0.0555  (2.3-8.8)
## 4 Orange      2019-~ Age g~ 0-4 years            0.0254  (0.3-4.7)
## 5 San Diego   2019-~ Age g~ 0-4 years            0.0180  (0.4-3.2)
## 6 Sutter      2019-~ Age g~ 0-4 years            0.00271 (0.1-0.4)
## # i abbreviated name: 1: '95% CONFIDENCE INTERVAL'
## # i 4 more variables: 'COUNTIES GROUPED' <chr>, COMMENT <chr>,
## #   CP_DATA_100 <dbl>, log_CP <dbl>
```

```r
# Use the tidy function from the broom package to obtain a clean summary of the model coefficients
tidy(asthma_lm)
```

```
## # A tibble: 4 x 5
##   term                estimate std.error statistic  p.value
##   <chr>                  <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)            0.778     0.156      4.97 1.82e- 6
## 2 AGE_GROUP18-64 years   1.53      0.164      9.29 1.84e-16
## 3 AGE_GROUP5-17 years    1.62      0.172      9.42 8.36e-17
## 4 AGE_GROUP65+ years     1.47      0.164      8.97 1.19e-15
```

```r
# Visualize the data with a scatter plot, coloring points by age group and adjusting themes for readabi
ggplot(asthma_pt3, aes(x=COUNTY, y=log_CP)) +
  geom_point(aes(col=AGE_GROUP)) +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, size=5)) +
  labs(
    x="Counties",
    y="Current Prevalence (LOG)",
    title="Current Prevalence of Asthma for Ages 0 to 65+ yrs. from 2019-2020"
  ) +
  facet_wrap(~AGE_GROUP, nrow=4)
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>
```
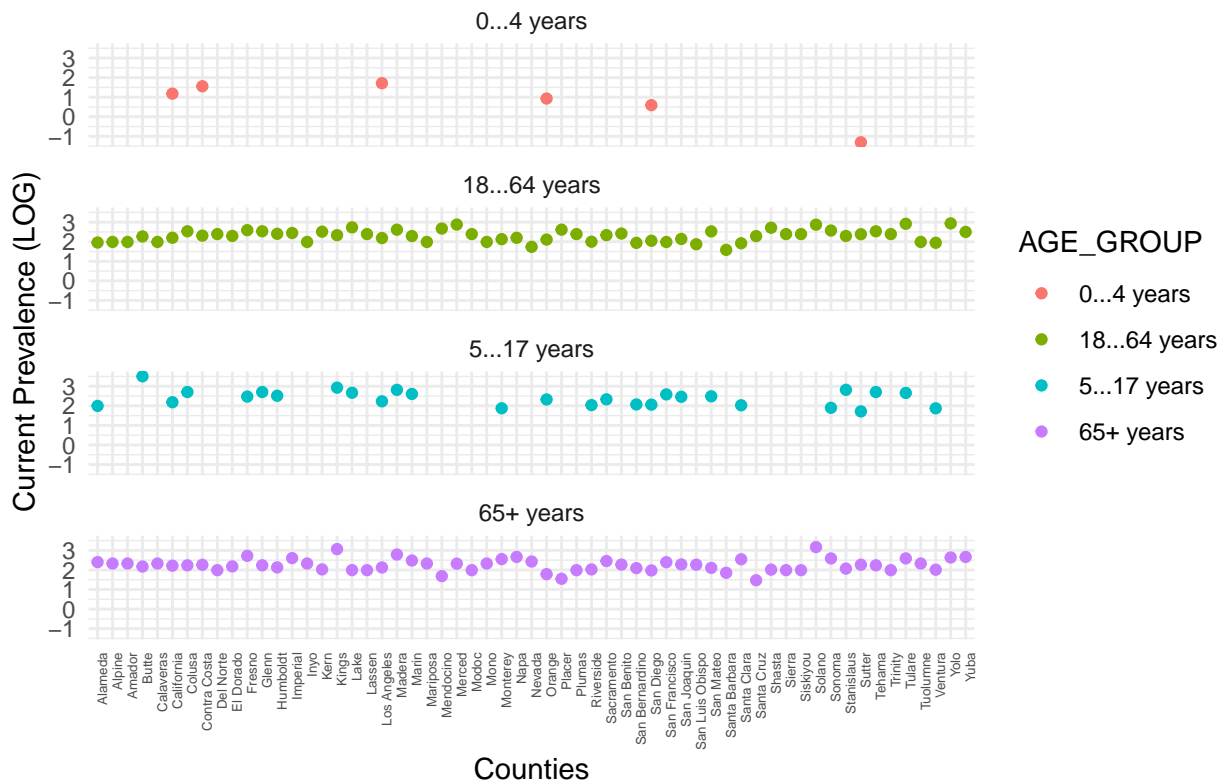
```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>
```
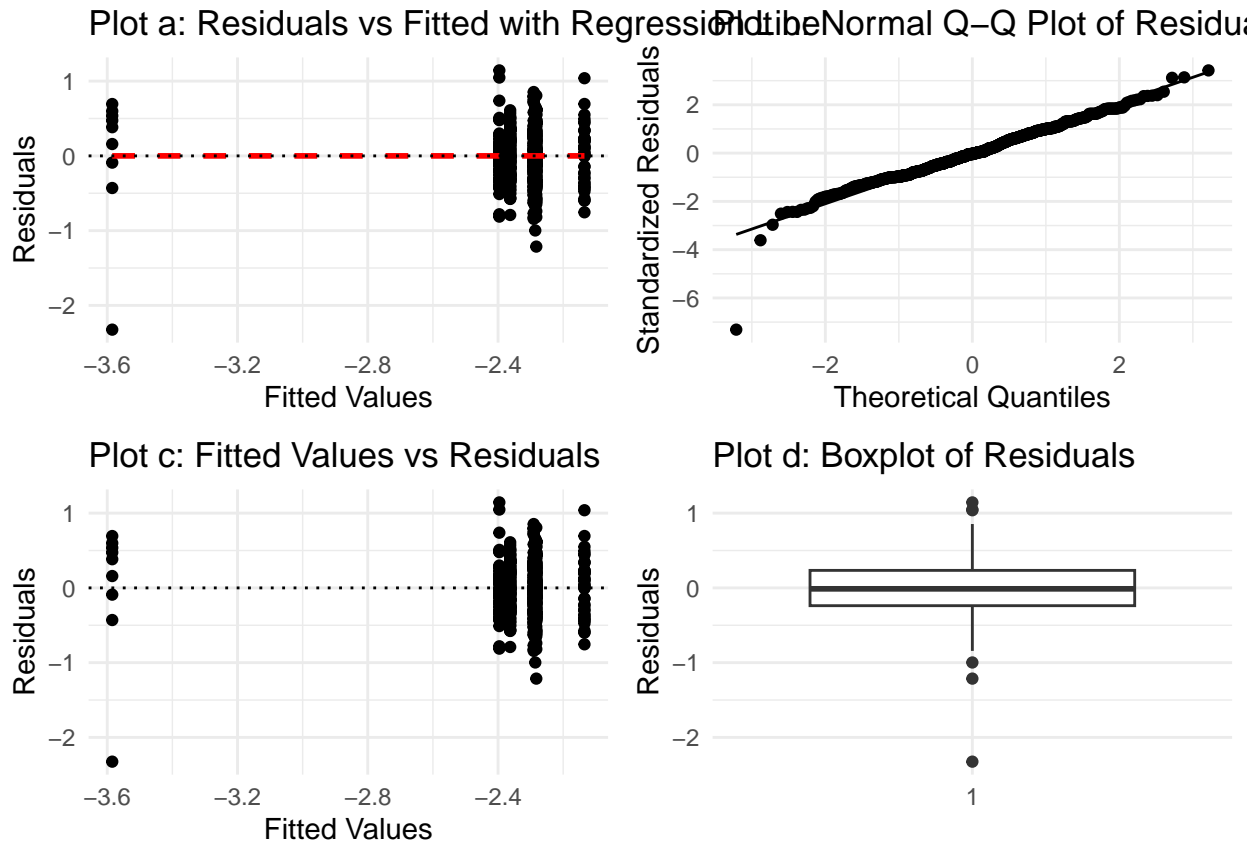
```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>
```

```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>
```



Current Prevalence of Asthma for Ages 0 to 65+ yrs. from 2019–2020

**21.** **[4 marks] Present your results in a clear summary.** This should include both a text summary and a table or figure with appropriate labeling. For example, if your outcome and predictor/exposure variables are both binary, this might be a 2x2 table. If your method was regression, you might present your regression line graphically.

```
#ggplot(asthma_lm,aes(sample=log_CP))+stat_qq_line()+stat_qq()+labs(x="Theoretical quantities",y="Resid
grid.arrange(plot_a, plot_b, plot_c, plot_d, nrow = 2)
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



```
ggplot(asthma_pt3, aes(x = COUNTY, y = log_CP, color = AGE_GROUP)) +
  geom_point(alpha = 0.6) +
  geom_smooth(method = "lm", se = FALSE, aes(group = AGE_GROUP), color = "black") +
  facet_wrap(~ AGE_GROUP, scales = "free_y") +
  labs(
    title = "Regression Analysis of Asthma Prevalence by Age Group",
    x = "County",
    y = "Log of Current Prevalence"
  ) +
  theme_minimal() +
  theme(
    axis.text.x = element_text(angle = 45, hjust = 1),
    plot.title = element_text(size = 14, face = "bold"),
    legend.position = "bottom"
  )
```

```
## 'geom_smooth()' using formula = 'y ~ x'

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>
```
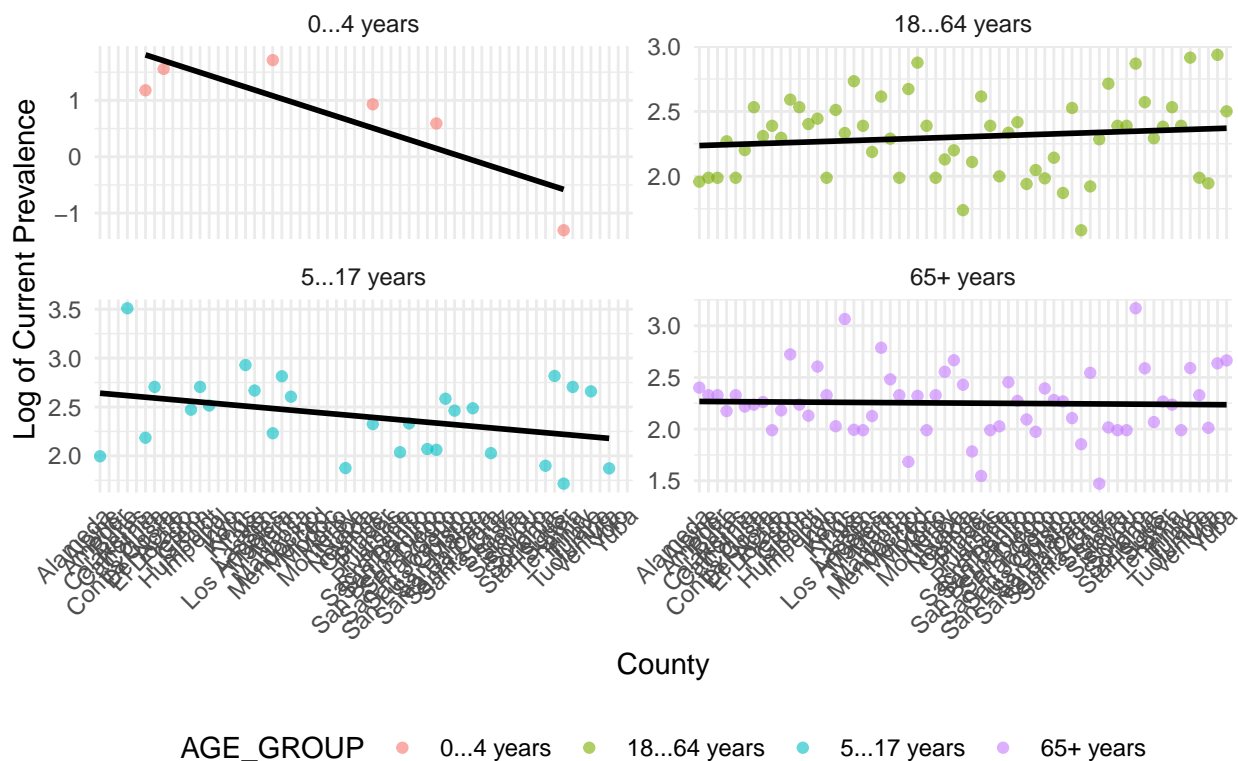
```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>
```

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <e2>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <80>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '0-4 years' in 'mbcsToSbcs': dot substituted for <93>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <e2>
```

```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '18-64 years' in 'mbcsToSbcs': dot substituted for <93>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <e2>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <80>


## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '5-17 years' in 'mbcsToSbcs': dot substituted for <93>
```



Regression Analysis of Asthma Prevalence by Age Group

```r
glance(asthma_lm)
```

```
## # A tibble: 1 x 12
##   r.squared adj.r.squared sigma statistic  p.value    df logLik   AIC   BIC
##       <dbl>         <dbl> <dbl>     <dbl>    <dbl> <dbl>  <dbl> <dbl> <dbl>
## 1     0.388         0.376 0.383      31.3 9.60e-16     3  -67.9  146.  161.
## # i 3 more variables: deviance <dbl>, df.residual <int>, nobs <int>
```

Our method was using a regression model to confirm our assumptions that are required for regression inference. We used diagnostic plots in regression analysis to help verify the model's assumptions and fit

60

by assessing linearity, variance consistency, normality of residuals, and identifying outliers through visual representations of residuals against predicted values, their distribution, and spread. Additionally, with a p-value of 9.6*10^-16 being less than 0.05, we can reject the null hypothesis. Therefore, we can assume that at least one predictor (AQI level/Current Prevalence or age group) significantly affects asthma rates (non-zero coefficients).

**22.** **[4 marks] Interpret your findings. Include a statement about the evidence, your conclusions, and the generalizability of your findings. Our analysis and conclusions depend on the quality of our study design and the methods of data collection. Any missteps or oversights during the data collection process could potentially change the outcome of what we are trying to find. Consider the methods used to collect the data you analyzed. Was there any potential issue in how the participants were selected/recruited, retained, or assessed that may have impacted the outcome of your analysis/visualization? Were there any potential biases that you might be concerned about? Were there factors that were not measured or considered that you think could be important to the interpretation of these data?**

Through our findings, we were able to conclude that there was at least one predictor, either AQI levels or age group, that significantly affects asthma rates (p-value=9.600795e-16) during 2019 to 2020. With further analysis and visualization using a scatter plot, we were able to conclude that age was a contributing factor to the prevalence of asthma. At each age group (5-17yrs., 18-64 yr., and 65+ yrs.) the p-values were significantly smaller (8.355944e-17,1.843716e-16,1.185769e-15) than the 0.05 threshold, therefore further revealing a relationship between age and asthma rates. With further analysis, our $R^2$ data (0.388) shows that our model predicts 38.8% relationship between asthma rates and age groups. The standard deviation of 0.383 shows that our data is moderately clustered around the mean. Our sample included 177 counties in California and this study can be generalized to populations similar to California, with further analysis focusing on a specific county during 2019 to 2020. Furthermore, if given the population size of each county, we could have analyzed a difference in asthma prevalence from year to year to see if there is any correlation. The data set that we used did not mention or reveal any information regarding race/ethnicity, income, education attainment, population size, and access to healthcare. We believe these additional variables are important to consider because, for example, populations with minimal access to healthcare will be under-reporting the prevalence of asthma. Additionally the data set did not mention how the data was collected, therefore, we were unable to conclude if there were any potential bias in sample collection. However, based on our graphs, we can assume that it follows within a simple random sample because it follows a normal distribution. Whether or not findings support rejecting the null hypothesis based on p-values how representative the sample is and to what populations results can be generalized potential biases or measurement errors that could affect validity of the results unmeasured confounders that might influence interpretation of the data.

**23. [1 mark] Create a statement of contribution. This is now common in journal articles for example the American Journal of Epidemiology provides the following instructions to authors: "Authorship credit should be based on criteria developed by the International Committee for Medical Journal Editors (ICMJE): 1) substantial contributions to conception and design, or acquisition of data, or analysis and interpretation of data; 2) drafting the article or reviewing it and, if appropriate, revising it critically for important intellectual content; 3) final approval of the version to be published. Authors should meet all conditions. In addition, each author must certify that he or she has participated sufficiently in the work to believe in its overall validity and to take public responsibility for appropriate portions of its content. Author names should be listed in ScholarOne and author contributions should be detailed in the cover letter (e.g., "Author A designed the study and directed its implementation, including quality assurance and control. Author B helped supervise the field activities and designed the study's analytic strategy. Author C helped conduct the literature review and prepare the Methods and the Discussion sections of the text.")."**

This project was completed through the collaborative efforts of each team member, who each contributed uniquely, but was finalized by Emma Collo and Adriel Vijuan.

Emma Collo led Part II of the project, focusing on the application of statistical distributions necessary for our analysis. Her role was crucial in ensuring that our analytical methods were aligned with the project's requirements and our group had thorough explanations in our responses. Emma also played a significant role in helping check in with our GSI, helping to clarify our approach and refine our findings for Part III. She also heavily contributed by analyzing our findings towards the latter half of Part III based on the statistical test used and kept communication between members on track. Adriel Vijuan provided a second-opinion on portions of the responses, particularly Part I and Part III, and helped refine our initial research question. He oversaw the framework structure of Part III's coding approach using regression inference and ensured that diagnostic plots were in place to allow the test to proceed properly.

Andres Carrillo Solis and Cleo Lin contributed to the early stages of the project, with Andres setting up the initial data framework and Cleo assisting in defining the research questions. Although Andres and Cleo had to leave the project early, their contributions in the initial discussions were foundational to the later portions of the work.

We would also like to thank our GSI for providing support by reviewing our project at crucial stages and offering feedback that helped keep our data analysis on track and valid per the requirements.

Thank you PH142!