# class 15

## Ramola Baviskar (PID A12228297)

## 3/8/2022

#Install datapasta

#Investigate Pertussis case numbers over time in the US

The CDC has tracked case numbers since the early 1920s. https://www.cdc.gov/pertussis/surv-reporting/cases-by-year.html

```
cdc <- data.frame(
                          Year = c(1922L,1923L,1924L,1925L,
                                   1926L,1927L,1928L,1929L,1930L,1931L,
                                   1932L,1933L,1934L,1935L,1936L,
                                   1937L,1938L,1939L,1940L,1941L,1942L,
                                   1943L,1944L,1945L,1946L,1947L,
                                   1948L,1949L,1950L,1951L,1952L,
                                   1953L,1954L,1955L,1956L,1957L,1958L,
                                   1959L,1960L,1961L,1962L,1963L,
                                   1964L,1965L,1966L,1967L,1968L,1969L,
                                   1970L,1971L,1972L,1973L,1974L,
                                   1975L,1976L,1977L,1978L,1979L,1980L,
                                   1981L,1982L,1983L,1984L,1985L,
                                   1986L,1987L,1988L,1989L,1990L,
                                   1991L,1992L,1993L,1994L,1995L,1996L,
                                   1997L,1998L,1999L,2000L,2001L,
                                   2002L,2003L,2004L,2005L,2006L,2007L,
                                   2008L,2009L,2010L,2011L,2012L,
                                   2013L,2014L,2015L,2016L,2017L,2018L,
                                   2019L),
         No..Reported.Pertussis.Cases = c(107473,164191,165418,152003,
                                   202210,181411,161799,197371,
                                   166914,172559,215343,179135,265269,
                                   180518,147237,214652,227319,103188,
                                   183866,222202,191383,191890,109873,
                                   133792,109860,156517,74715,69479,
                                   120718,68687,45030,37129,60886,
                                   62786,31732,28295,32148,40005,
                                   14809,11468,17749,17135,13005,6799,
                                   7717,9718,4810,3285,4249,3036,
                                   3287,1759,2402,1738,1010,2177,2063,
                                   1623,1730,1248,1895,2463,2276,
                                   3589,4195,2823,3450,4157,4570,
                                   2719,4083,6586,4617,5137,7796,6564,
                                   7405,7298,7867,7580,9771,11647,
                                   25827,25616,15632,10454,13278,
```

```
                                        16858,27550,18719,48277,28639,32971,
                                        20762,17972,18975,15609,18617)
       )
```

#Now use ggplot.

```
library(ggplot2)
library(tidyverse)
```

```
## -- Attaching packages ----------------------------------------- tidyverse 1.3.1 --

## v tibble   3.1.6      v dplyr    1.0.8
## v tidyr    1.2.0      v stringr 1.4.0
## v readr    2.1.2      v forcats 0.5.1
## v purrr    0.3.4

## -- Conflicts ------------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

Q1. Q2.

```
Pertussis_linegraph <- ggplot(cdc) +
  aes(Year, No..Reported.Pertussis.Cases) +
  geom_point() +
  geom_line()
```

```
Pertussis_linegraph +
  geom_vline(xintercept = 1946, color = "red", size = 1, linetype = "dashed") +
  geom_vline(xintercept = 1996, color = "blue", size= 1, linetype = "dashed")
```

Q3. Rates of pertussis increased after the aP vaccine. Possible reasons for this are: vaccine hesitancy, evolution of B. pertussis, increased testing, and a decreasing immunity among those vaccinated with the aP vaccine rather than the wP vaccine.

```
library(jsonlite)
```

```
##
## Attaching package: 'jsonlite'
```

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

#Exploring CMI-PDB data We'll use the **jsonlite** package to read from the CMI-PB database API directly.

```
url <- "https://www.cmi-pb.org/api/subject"

subject <- read_json(url, simplifyVector = TRUE)
head(subject, 3)
```

```
##   subject_id infancy_vac biological_sex              ethnicity  race
## 1          1          wP        Female Not Hispanic or Latino White
## 2          2          wP        Female Not Hispanic or Latino White
```

```
## 3             3          wP          Female                      Unknown White
##   year_of_birth date_of_boost   study_name
## 1    1986-01-01    2016-09-12 2020_dataset
## 2    1968-01-01    2019-01-28 2020_dataset
## 3    1983-01-01    2016-10-10 2020_dataset
```

```
table(subject$infancy_vac)
```

```
##
## aP wP
## 47 49
```

```
nrow(subject)
```

```
## [1] 96
```

Q4. ap: 47 wP: 49

```
table(subject$biological_sex)
```

```
##
## Female   Male
##     66     30
```

Q5. Female: 66 Male: 30

```
table(subject$biological_sex, subject$race)
```

```
##
##         American Indian/Alaska Native Asian Black or African American
##   Female                            0    18                         2
##   Male                             1     9                         0
##
##         More Than One Race Native Hawaiian or Other Pacific Islander
##   Female                  8                                        1
##   Male                    2                                        1
##
##         Unknown or Not Reported White
##   Female                      10    27
##   Male                         4    13
```

Q6. Female American Indian/Alaska Native: 0 Female Asian: 18 Female Black/African American: 2 Female More Than One Race: 8 Female Native Hawaiian/Other Pac. Islander: 1 Female Unknown/Not Reported: 10 Female White: 27 Male American Indian/Alaska Native: 1 Male Asian: 9 Male Black/African American: 0 Male More Than One Race: 2 Male Native Hawaiian/Other Pac. Islander: 1 Male Unknown/Not Reported: 4 Male White: 13

```
library(lubridate)
```

```
## 
## Attaching package: 'lubridate'

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

Q7. $ Q8. optional

#Join datasets.

```
specimen <- read_json("https://www.cmi-pb.org/api/specimen", simplifyVector = TRUE)
titer <- read_json("https://www.cmi-pb.org/api/ab_titer", simplifyVector = TRUE)
```

Take a quick look.

```
head(specimen, 3)
```

```
##   specimen_id subject_id actual_day_relative_to_boost
## 1           1          1                           -3
## 2           2          1                          736
## 3           3          1                            1
##   planned_day_relative_to_boost specimen_type visit
## 1                             0         Blood     1
## 2                           736         Blood    10
## 3                             1         Blood     2
```

I need to use inner_join() here.

Q9.

```
library(dplyr)
library(tidyverse)
```

Q9.

```
meta <- inner_join(specimen, subject)
```

```
## Joining, by = "subject_id"
```

```
dim(meta)
```

```
## [1] 729  13
```

```
head(meta)
```

```
##   specimen_id subject_id actual_day_relative_to_boost
## 1           1          1                           -3
## 2           2          1                          736
## 3           3          1                            1
## 4           4          1                            3
## 5           5          1                            7
## 6           6          1                           11
##   planned_day_relative_to_boost specimen_type visit infancy_vac biological_sex
## 1                             0         Blood     1          wP         Female
## 2                           736         Blood    10          wP         Female
## 3                             1         Blood     2          wP         Female
## 4                             3         Blood     3          wP         Female
## 5                             7         Blood     4          wP         Female
## 6                            14         Blood     5          wP         Female
##               ethnicity  race year_of_birth date_of_boost    study_name
## 1 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 2 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 3 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 4 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 5 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 6 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
```

Q10.

```
abdata <- inner_join(titer, meta)
```

```
## Joining, by = "specimen_id"
```

```
dim(abdata)
```

```
## [1] 32675    19
```

```
table(abdata$isotype)
```

```
##
##  IgE  IgG IgG1 IgG2 IgG3 IgG4
## 6698 1413 6141 6141 6141 6141
```

Q11. IgE: 6698 IgG:1413 IgG1: 6141 IgG2:6141 IgG3:6141 IgG4: 6141

```
table(abdata$visit)
```

```
##
##    1    2    3    4    5    6    7    8
## 5795 4640 4640 4640 4640 4320 3920   80
```

Q12. There are very vew visit 8 specimens compared to other visits. It's likely unfinished.
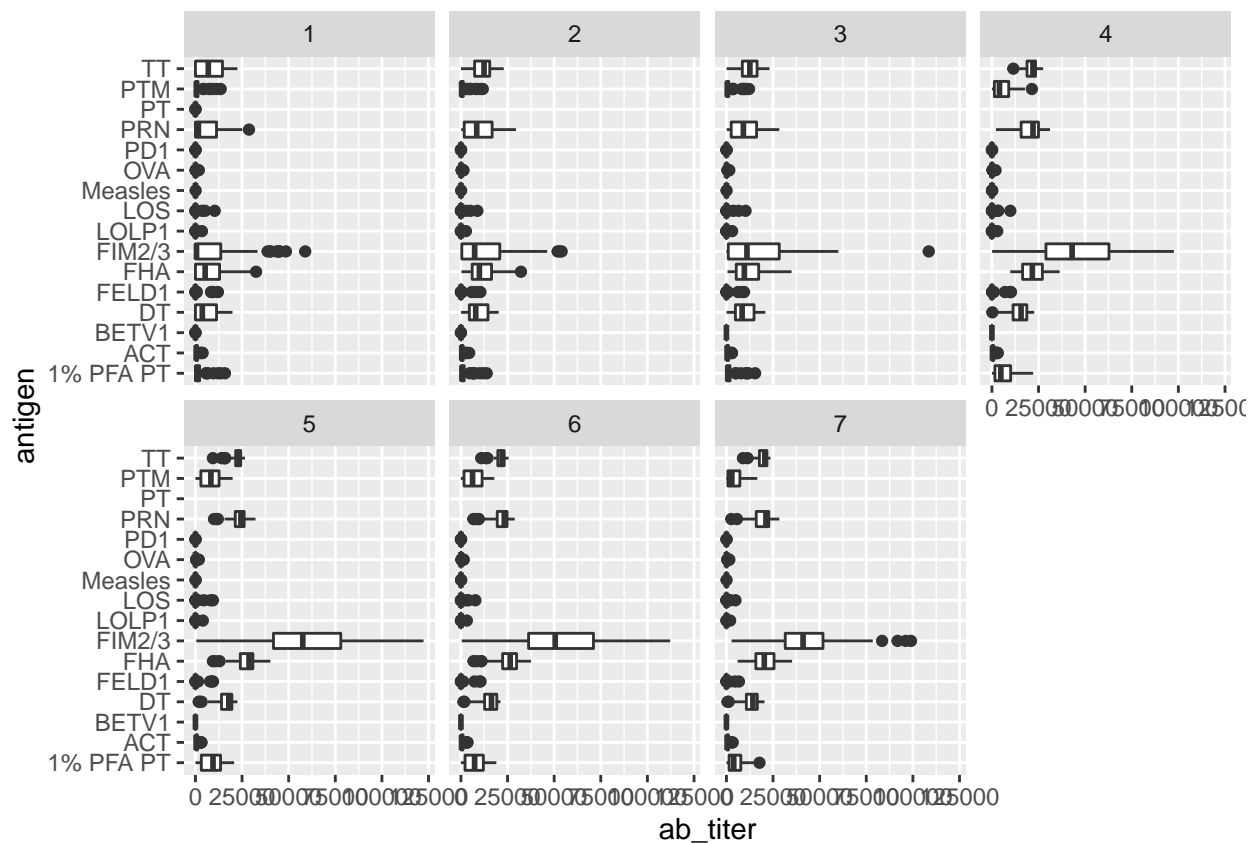
Now we'll exclude visit 8 because it is unfinished.

```
ig1 <- abdata %>% filter(isotype == "IgG1", visit!=8)
head(ig1)
```

```
##   specimen_id isotype is_antigen_specific antigen    ab_titer  unit
## 1           1    IgG1                TRUE     ACT 274.355068 IU/ML
## 2           1    IgG1                TRUE     LOS  10.974026 IU/ML
## 3           1    IgG1                TRUE   FELD1   1.448796 IU/ML
## 4           1    IgG1                TRUE   BETV1   0.100000 IU/ML
## 5           1    IgG1                TRUE   LOLP1   0.100000 IU/ML
## 6           1    IgG1                TRUE Measles  36.277417 IU/ML
##   lower_limit_of_detection subject_id actual_day_relative_to_boost
## 1                 3.848750          1                           -3
## 2                 4.357917          1                           -3
## 3                 2.699944          1                           -3
## 4                 1.734784          1                           -3
## 5                 2.550606          1                           -3
## 6                 4.438966          1                           -3
##   planned_day_relative_to_boost specimen_type visit infancy_vac biological_sex
## 1                             0         Blood     1          wP         Female
## 2                             0         Blood     1          wP         Female
## 3                             0         Blood     1          wP         Female
## 4                             0         Blood     1          wP         Female
## 5                             0         Blood     1          wP         Female
## 6                             0         Blood     1          wP         Female
##                 ethnicity  race year_of_birth date_of_boost   study_name
## 1 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 2 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 3 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 4 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 5 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
## 6 Not Hispanic or Latino White    1986-01-01    2016-09-12 2020_dataset
```
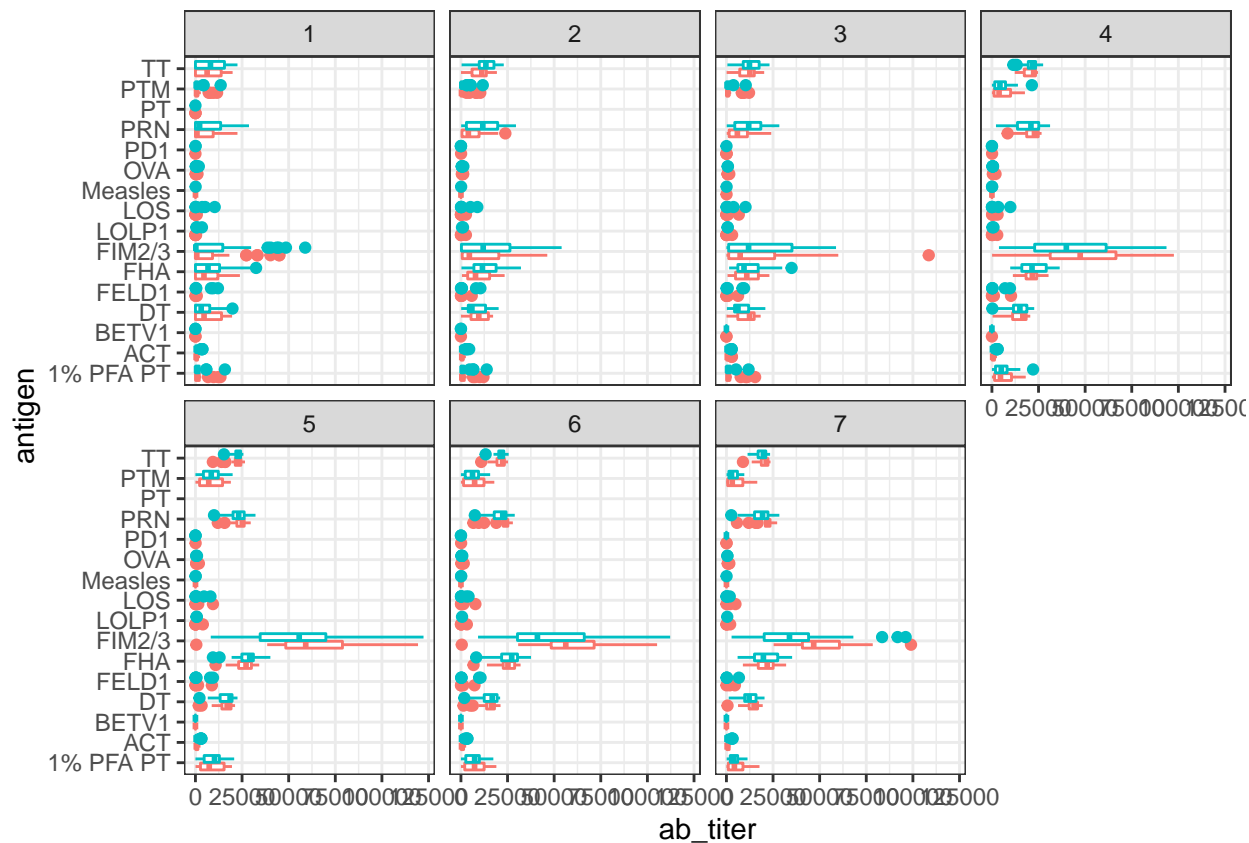
Q13.

```
ggplot(ig1) +
  aes(ab_titer, antigen) +
  geom_boxplot() +
  facet_wrap(vars(visit), nrow=2)
```
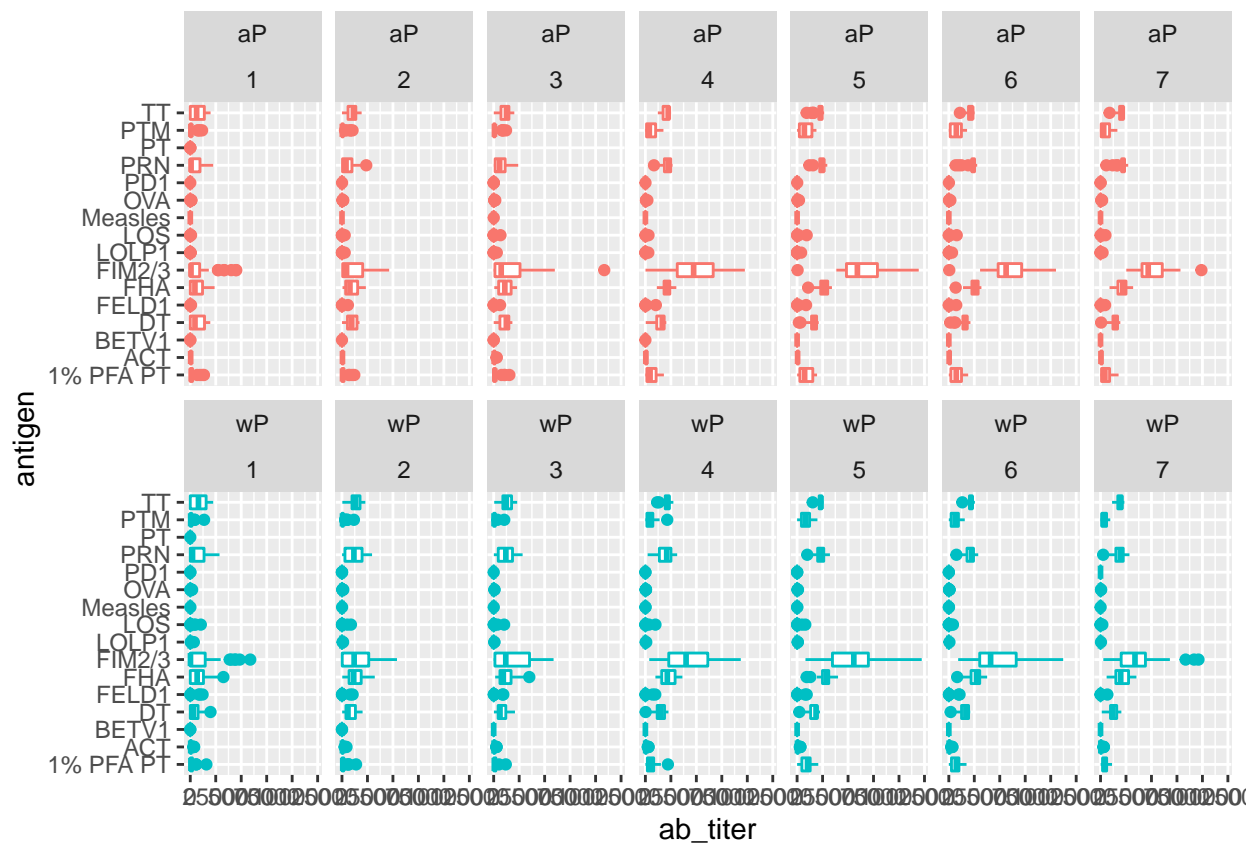
Q14. The FIM2/3 antigen has shifted. This is likely because antibodies have specifically recognized it.

```
ggplot(ig1) +
  aes(ab_titer, antigen, col=infancy_vac ) +
  geom_boxplot(show.legend = FALSE) +
  facet_wrap(vars(visit), nrow=2) +
  theme_bw()
```
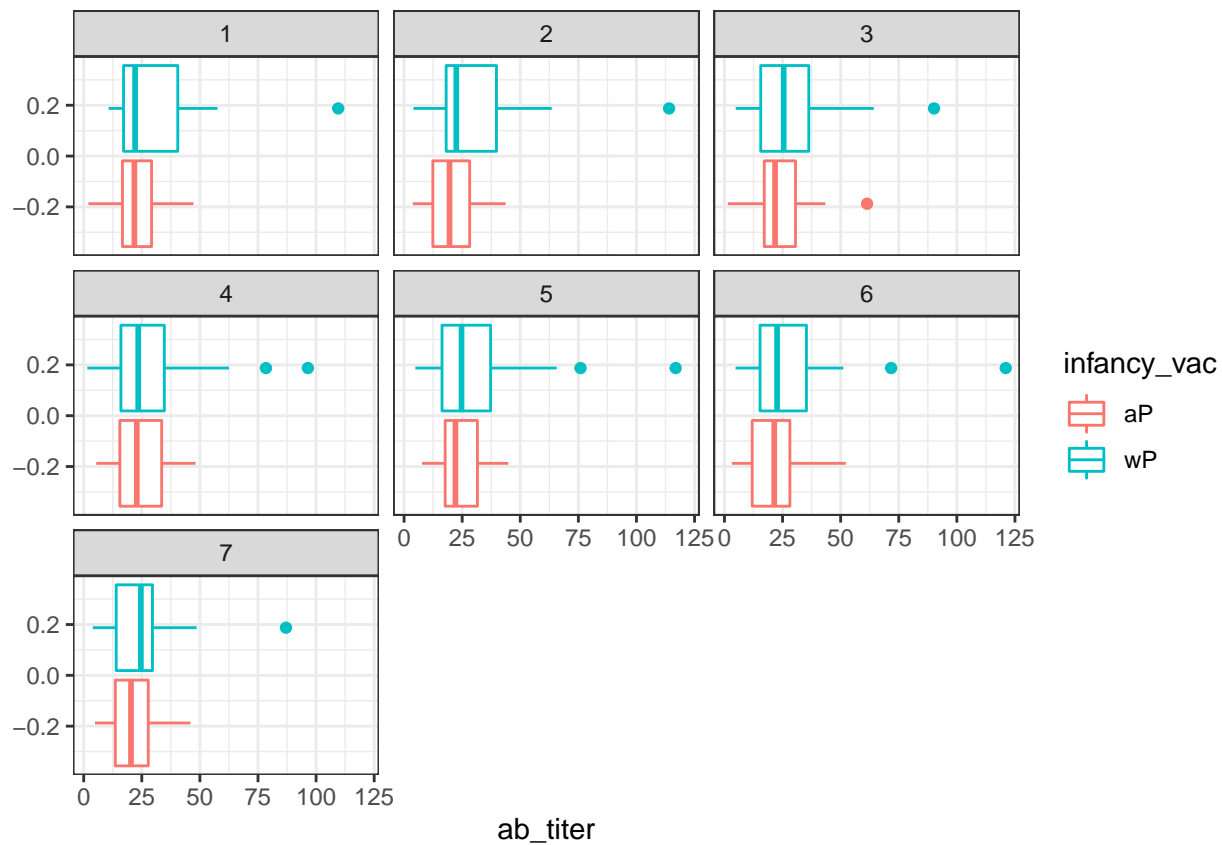
```
ggplot(ig1) +
  aes(ab_titer, antigen, col=infancy_vac ) +
  geom_boxplot(show.legend = FALSE) +
  facet_wrap(vars(infancy_vac, visit), nrow=2)
```
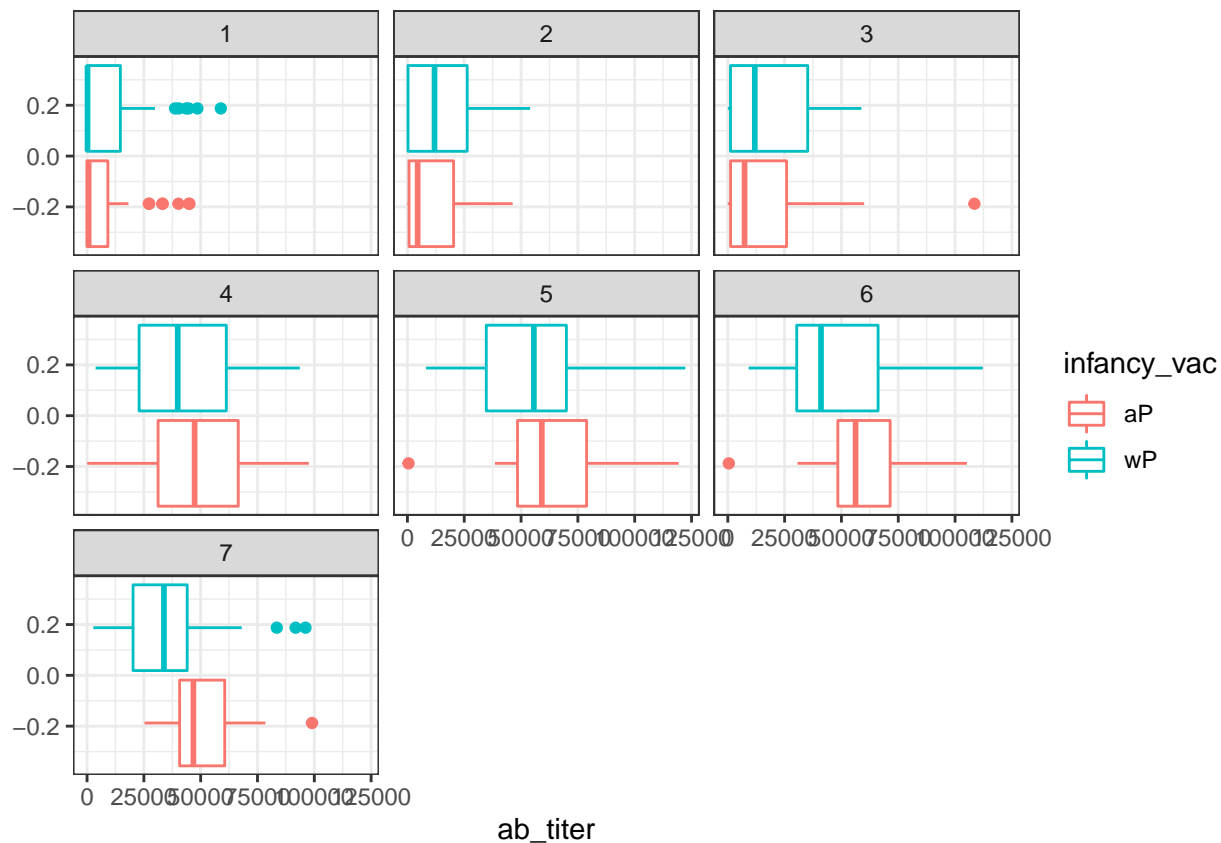
9

Q15.

```
filter(ig1, antigen=="Measles") %>%
  ggplot() +
  aes(ab_titer, col=infancy_vac) +
  geom_boxplot(show.legend = TRUE) +
  facet_wrap(vars(visit)) +
  theme_bw()
```

```
filter(ig1, antigen=="FIM2/3") %>%
  ggplot() +
  aes(ab_titer, col=infancy_vac) +
  geom_boxplot(show.legend = TRUE) +
  facet_wrap(vars(visit)) +
  theme_bw()
```

ab_titer

>Q16. The measles course is remarkable steady. It scarcely changes at all through the 8 visits. The FIM2/3 data, however, shows quite a lot of change. In both aP and wpP trials, it rises pretty consistently until visit 5, after which there is a slight decline.
>Q17. No.

# Pull RNA-Seq data from the CMI-PB database.

We can use the CMI-PB API to pull obtain time-course RNA-Seq results for wP and aP subjects (i.e. patients).

```
url <- "https://www.cmi-pb.org/api/v2/rnaseq?versioned_ensembl_gene_id=eq.ENSG00000211896.7"
rna <- read_json(url, simplifyVector = TRUE)
```
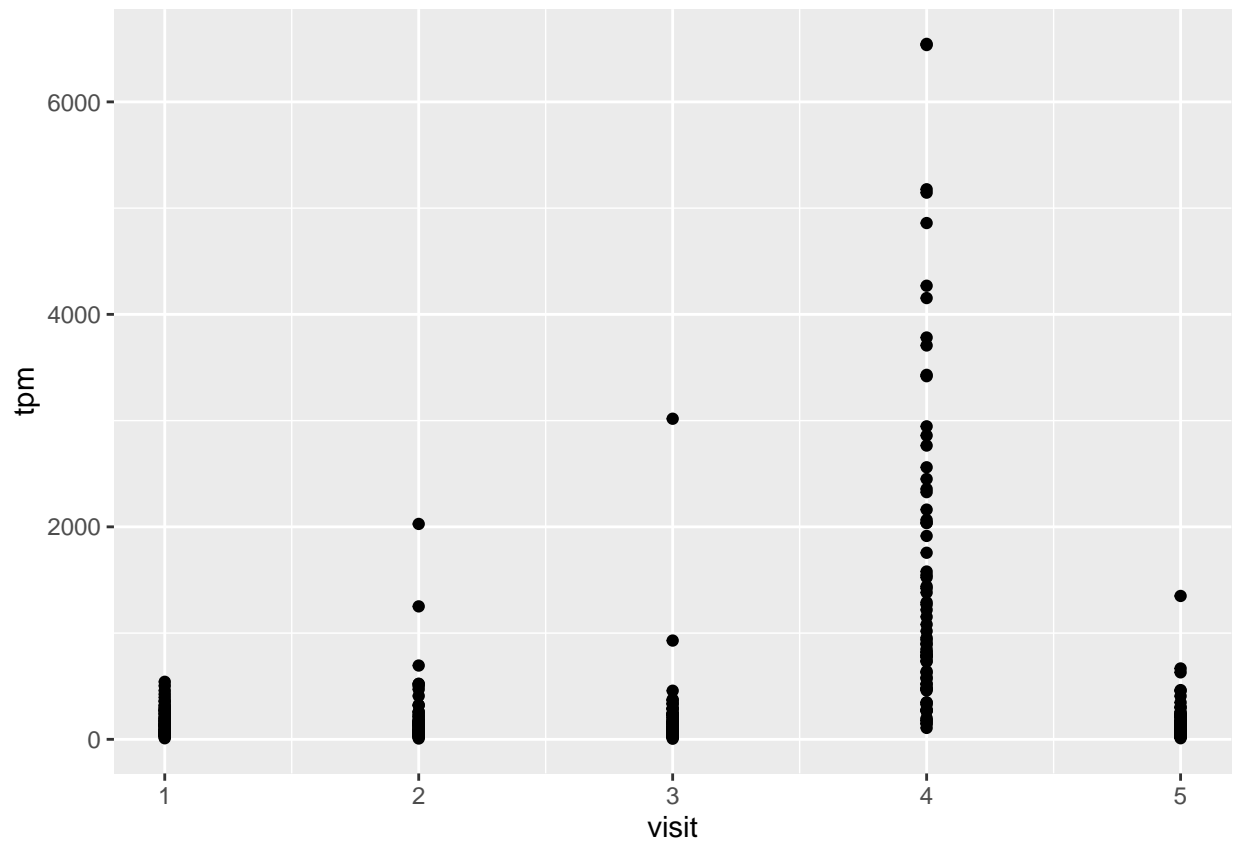
```
ssrna <- inner_join(rna, meta)
```

```
## Joining, by = "specimen_id"
```
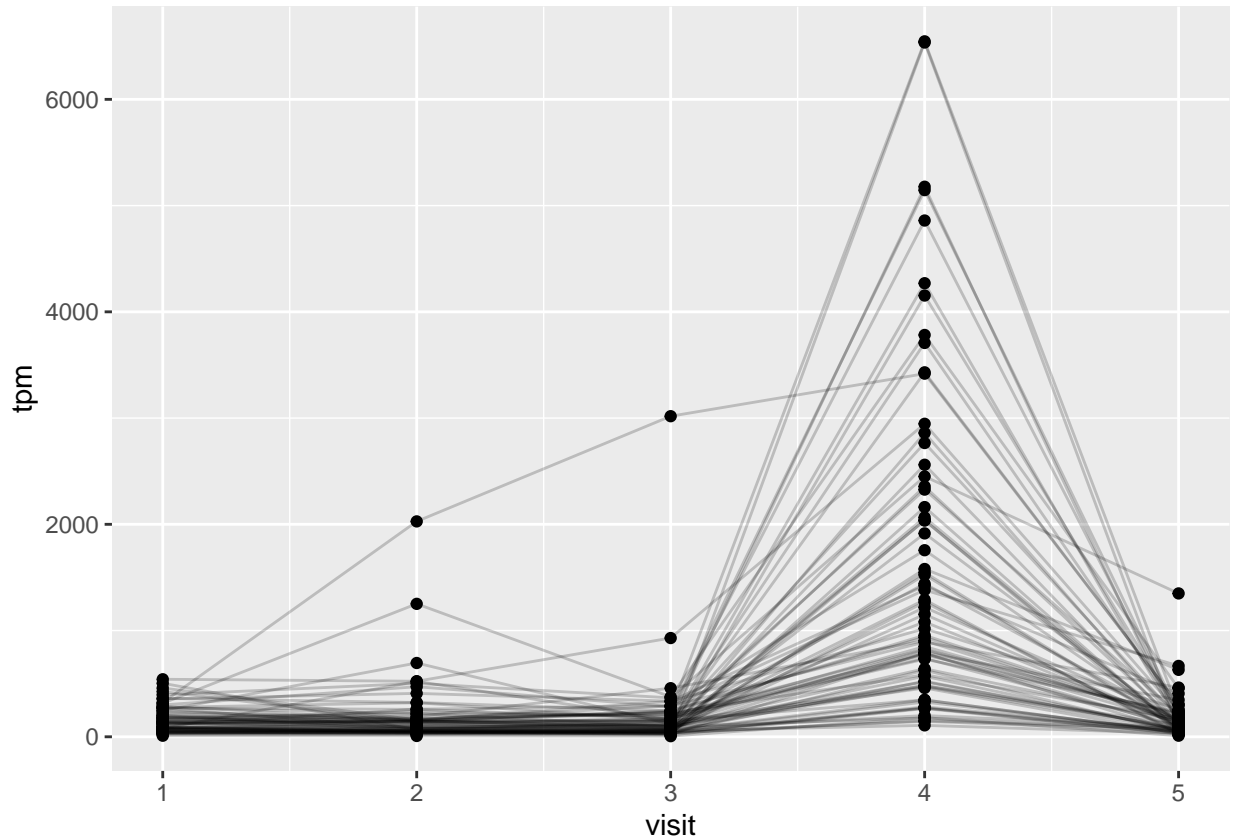
```
dim(ssrna)
```

```
## [1] 360  16
```

Q18.

12

```
ggplot(ssrna) +
  aes(visit, tpm) +
  geom_point()
```



```
ggplot(ssrna) +
  aes(visit, tpm, group=subject_id) +
  geom_point() +
  geom_line(alpha=0.2)
```

Q19. It's at its maximum at around visit 4. Q20. It sort of matches. The AB Titer data suggested a peak at around visit 5 while the gene peaks at visit 4. The gene expression leads to the creation of antibodies; once a sufficient quantity of the antibody has been manufactured, the cell expression drops off. At this point, many antibodies are present (peaking at visit 5) and persist for some time.

```
ggplot(ssrna) +
  aes(tpm, col=infancy_vac) +
  geom_boxplot() +
  facet_wrap(vars(visit))
```