# Class 19: Pertussis and the CMI-PB project

Jiachen Fan (A17662703)

## 1. Investigating pertussis cases by year

We can view this data on the CDC website here: https://www.cdc.gov/pertussis/surv-reporting/cases-by-year.html > Q1. With the help of the R "addin" package datapasta assign the CDC pertussis case number data to a data frame called cdc and use ggplot to make a plot of cases numbers over time.

```
# install.packages("datapasta")
```

```
library(tidyverse)
```

Warning: package 'tidyverse' was built under R version 4.3.2

Warning: package 'readr' was built under R version 4.3.2

Warning: package 'forcats' was built under R version 4.3.2

Warning: package 'lubridate' was built under R version 4.3.2

```
-- Attaching core tidyverse packages ----------------------- tidyverse 2.0.0 --
v dplyr     1.1.3     v readr     2.1.4
v forcats   1.0.0     v stringr   1.5.0
v ggplot2   3.4.4     v tibble    3.2.1
v lubridate 1.9.3     v tidyr     1.3.0
v purrr     1.0.2
-- Conflicts ------------------------------------------ tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to becom
```
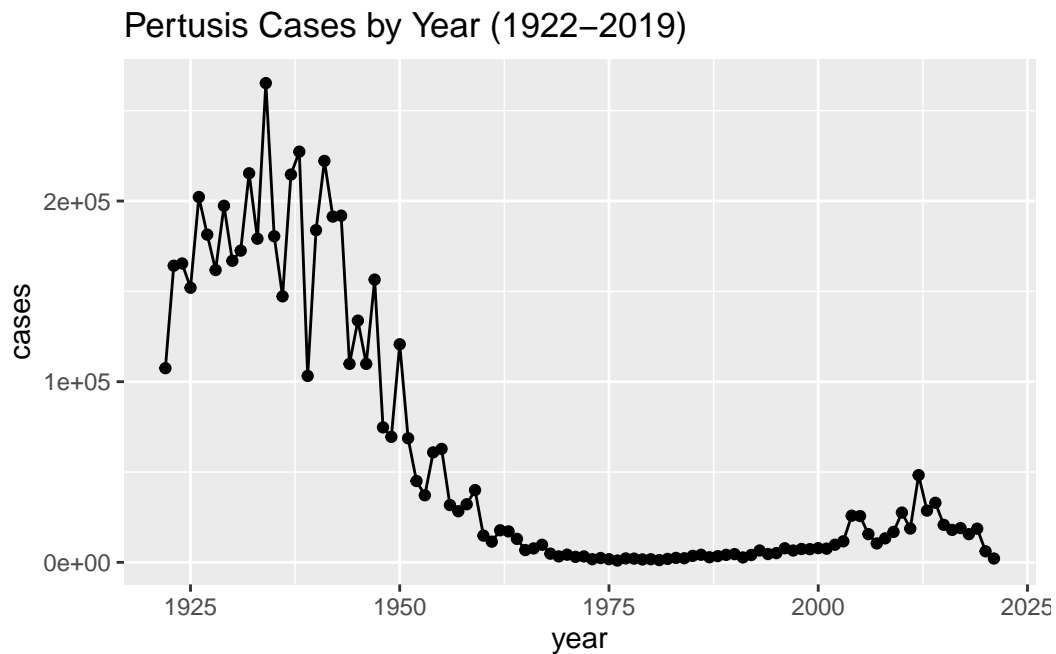
```
library(ggplot2)

ggplot(cdc) +
  aes(x= year,y= cases) +
  geom_point() +
  geom_line() +
  labs(title = "Pertusis Cases by Year (1922-2019)")
```
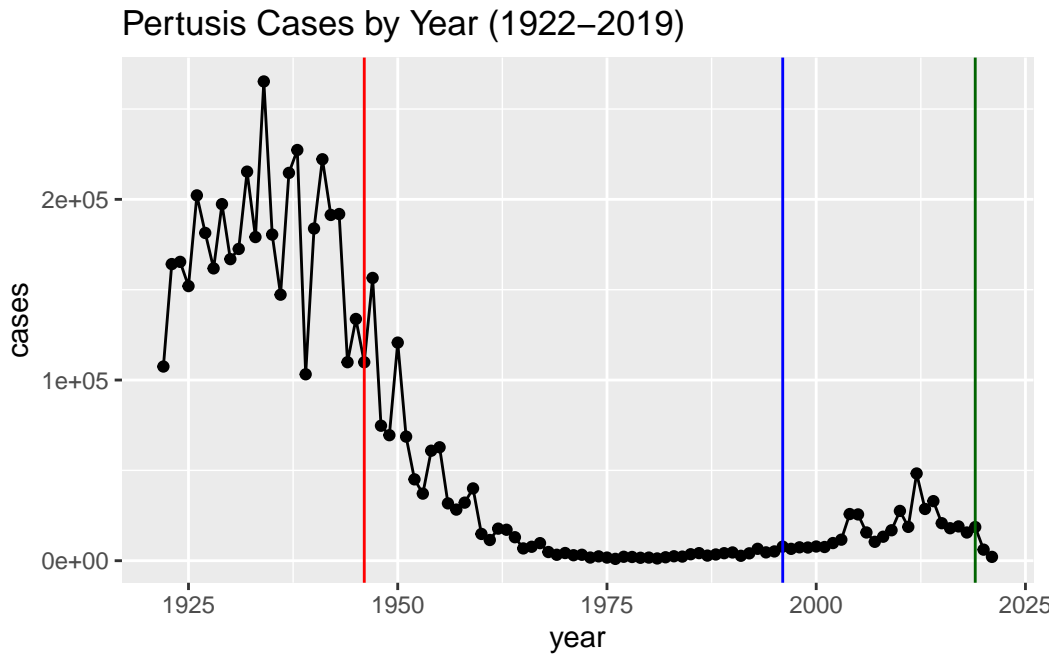
Pertusis Cases by Year (1922–2019)



## 2. A tale of two vaccines (wP & aP)

Q2. Using the ggplot geom_vline() function add lines to your previous plot for
the 1946 introduction of the wP vaccine and the 1996 switch to aP vaccine (see
example in the hint below). What do you notice?

```
ggplot(cdc) +
  aes(x= year,y= cases) +
  geom_point() +
  geom_line() +
  geom_vline(xintercept=1946, color ="red") +
  geom_vline(xintercept=1996, color ="blue") +
```

```
geom_vline(xintercept=2019, color ="darkgreen")+
labs(title = "Pertusis Cases by Year (1922-2019)")
```



Pertusis Cases by Year (1922–2019)

Reported cases reduced significantly after introduction of wP vaccine.

> Q3. Describe what happened after the introduction of the aP vaccine? Do you have a possible explanation for the observed trend?

1) more sensitive PCR-based testing, 2) vaccination hesitancy 3) bacterial evolution (escape from vaccine immunity) 4) waning of immunity in adolescents originally primed as infants with the newer aP vaccine as compared to the older wP vaccine.

Additional points for discussion: How are vaccines currently approved?

Typically we first examine 'Correlates of protection' which are things that can be measured within weeks or months after vaccination, and which are thought to correlate with increased protection from disease. For the aP vaccine this was an induction of antibodies against pertussis toxin (PT) in infants at equivalent levels to those induced by the wP vaccine. The aP vaccines also had less side effects (reduction of sore arms, fever and pain). Testing for protection induced by a new vaccine requires a lot of people exposed to the pathogen (like in a pandemic). It is impossible to discover a effect 10 years post vaccination in the current trial system. It is unclear what differentiates people that have been primed with aP vs. wP long

term. The CMI-PB project is an attempt to make data on this question open and examinable by all.

## 3. Exploring CMI-PB data

```
library(jsonlite)
```

Attaching package: 'jsonlite'

The following object is masked from 'package:purrr':

    flatten

```
subject <- read_json("https://www.cmi-pb.org/api/subject", 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     dataset
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
```

Q4. How many aP and wP infancy vaccinated subjects are in the dataset?

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

```
aP wP
60 58
```

There are 60 aP and 58 wP in the dataset.

Q5. How many Male and Female subjects/patients are in the dataset?

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

```
Female    Male
    79      39
```

There are 79 Female and 39 Male patients in the dataset.

> Q6. What is the breakdown of race and biological sex (e.g. number of Asian females, White males etc...)?

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

```
         American Indian/Alaska Native Asian Black or African American
Female                              0    21                          2
Male                                1    11                          0

         More Than One Race Native Hawaiian or Other Pacific Islander
Female                    9                                         1
Male                      2                                         1

         Unknown or Not Reported White
Female                        11    35
Male                           4    20
```

> Q. Make a histogram of the subject age distribution and facet by infancy_vac

```r
library(lubridate)
```
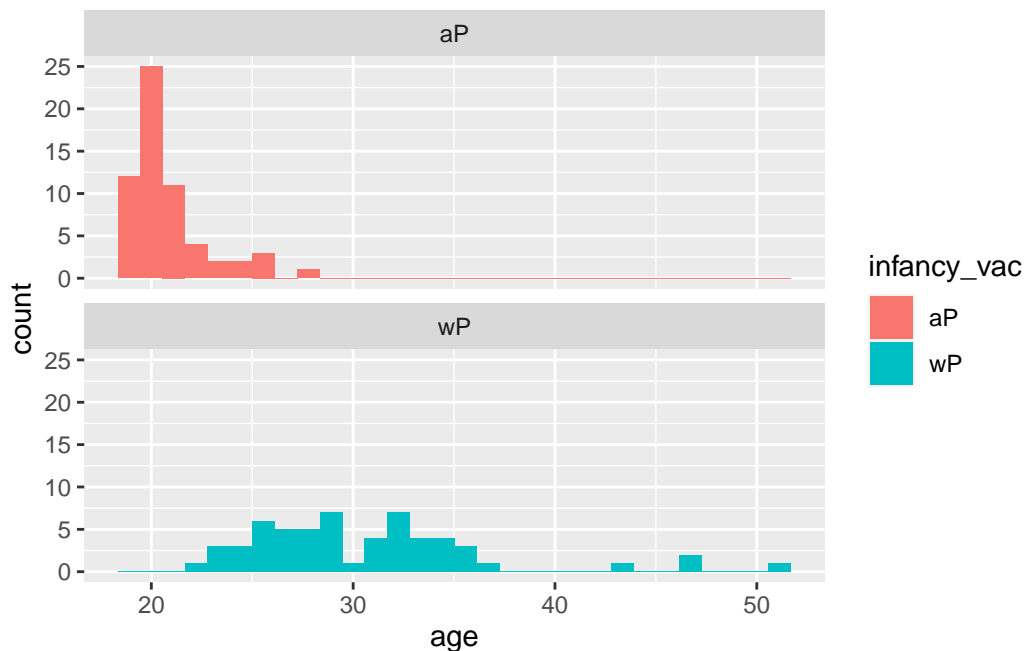
What is today's date

```r
today()
```

```
[1] "2023-12-11"
```

```r
subject$age <- time_length(ymd(subject$date_of_boost)-ymd(subject$year_of_birth), "years")
```

```
ggplot(subject)+
  aes(age,
      fill = infancy_vac)+
  facet_wrap(vars(infancy_vac),nrow = 2)+
  geom_histogram()
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



Q7. Using this approach determine (i) the average age of wP individuals, (ii) the average age of aP individuals; and (iii) are they significantly different?

```
subject$age_today <- time_length(today()-ymd(subject$year_of_birth),"years")
round(mean(subject$age_today[subject$infancy_vac=="wP"]))
```

[1] 36

```
round(mean(subject$age_today[subject$infancy_vac=="aP"]))
```

[1] 26

```r
t.test(subject$age_today[subject$infancy_vac=="wP"],subject$age_today[subject$infancy_vac=
```

```
[1] 6.813505e-19
```

i) 36; ii) 26; iii) Yes

```r
table(subject$dataset)
```

```
2020_dataset 2021_dataset 2022_dataset
          60           36           22
```

Joining multiple tables

```r
specimen <- read_json("https://www.cmi-pb.org/api/specimen", simplifyVector = TRUE)
titer <- read_json("https://www.cmi-pb.org/api/v4/plasma_ab_titer", simplifyVector = TRUE)
```

```r
head(specimen)
```

```
  specimen_id subject_id actual_day_relative_to_boost
1           1          1                           -3
2           2          1                            1
3           3          1                            3
4           4          1                            7
5           5          1                           11
6           6          1                           32
  planned_day_relative_to_boost specimen_type visit
1                             0         Blood     1
2                             1         Blood     2
3                             3         Blood     3
4                             7         Blood     4
5                            14         Blood     5
6                            30         Blood     6
```

Q9. Complete the code to join specimen and subject tables to make a new merged data frame containing all specimen records along with their associated subject details:

```r
library(dplyr)

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

Joining with `by = join_by(subject_id)`

```r
dim(meta)
```

[1] 939  15

```r
head(meta)
```

```
  specimen_id subject_id actual_day_relative_to_boost
1           1          1                           -3
2           2          1                            1
3           3          1                            3
4           4          1                            7
5           5          1                           11
6           6          1                           32
  planned_day_relative_to_boost specimen_type visit infancy_vac biological_sex
1                             0         Blood     1          wP         Female
2                             1         Blood     2          wP         Female
3                             3         Blood     3          wP         Female
4                             7         Blood     4          wP         Female
5                            14         Blood     5          wP         Female
6                            30         Blood     6          wP         Female
              ethnicity  race year_of_birth date_of_boost       dataset
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
       age age_today
1 30.69678  37.94114
2 30.69678  37.94114
3 30.69678  37.94114
4 30.69678  37.94114
5 30.69678  37.94114
6 30.69678  37.94114
```

Q10. Now using the same procedure join meta with titer data so we can further analyze this data in terms of time of visit aP/wP, male/female etc.

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

```
Joining with `by = join_by(specimen_id)`
```

```
dim(abdata)
```

```
[1] 41810    22
```

```
head(abdata)
```

```
  specimen_id isotype is_antigen_specific antigen        MFI MFI_normalised
1           1     IgE               FALSE   Total 1110.21154       2.493425
2           1     IgE               FALSE   Total 2708.91616       2.493425
3           1     IgG                TRUE      PT   68.56614       3.736992
4           1     IgG                TRUE     PRN  332.12718       2.602350
5           1     IgG                TRUE     FHA 1887.12263      34.050956
6           1     IgE                TRUE     ACT    0.10000       1.000000
  unit lower_limit_of_detection subject_id actual_day_relative_to_boost
1 UG/ML                 2.096133          1                           -3
2 IU/ML                29.170000          1                           -3
3 IU/ML                 0.530000          1                           -3
4 IU/ML                 6.205949          1                           -3
5 IU/ML                 4.679535          1                           -3
6 IU/ML                 2.816431          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      dataset
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
       age age_today
1 30.69678  37.94114
2 30.69678  37.94114
3 30.69678  37.94114
4 30.69678  37.94114
5 30.69678  37.94114
6 30.69678  37.94114
```

Q11. How many specimens (i.e. entries in abdata) do we have for each isotype?

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

```
 IgE  IgG IgG1 IgG2 IgG3 IgG4
6698 3240 7968 7968 7968 7968
```

Q12. What are the different $dataset values in abdata and what do you notice about the number of rows for the most "recent" dataset?

```
table(abdata$dataset)
```

```
2020_dataset 2021_dataset 2022_dataset
       31520         8085         2205
```

## 4. Examine IgG Ab titer levels

```
igg <- abdata %>% filter(isotype == "IgG")
head(igg)
```

```
  specimen_id isotype is_antigen_specific antigen       MFI MFI_normalised
1           1     IgG                TRUE      PT   68.56614       3.736992
2           1     IgG                TRUE     PRN  332.12718       2.602350
3           1     IgG                TRUE     FHA 1887.12263      34.050956
4          19     IgG                TRUE      PT   20.11607       1.096366
5          19     IgG                TRUE     PRN  976.67419       7.652635
6          19     IgG                TRUE     FHA   60.76626       1.096457
```

```
    unit lower_limit_of_detection subject_id actual_day_relative_to_boost
1 IU/ML                  0.530000          1                            -3
2 IU/ML                  6.205949          1                            -3
3 IU/ML                  4.679535          1                            -3
4 IU/ML                  0.530000          3                            -3
5 IU/ML                  6.205949          3                            -3
6 IU/ML                  4.679535          3                            -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       dataset
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                Unknown White    1983-01-01    2016-10-10 2020_dataset
5                Unknown White    1983-01-01    2016-10-10 2020_dataset
6                Unknown White    1983-01-01    2016-10-10 2020_dataset
       age age_today
1 30.69678  37.94114
2 30.69678  37.94114
3 30.69678  37.94114
4 33.77413  40.94182
5 33.77413  40.94182
6 33.77413  40.94182
```
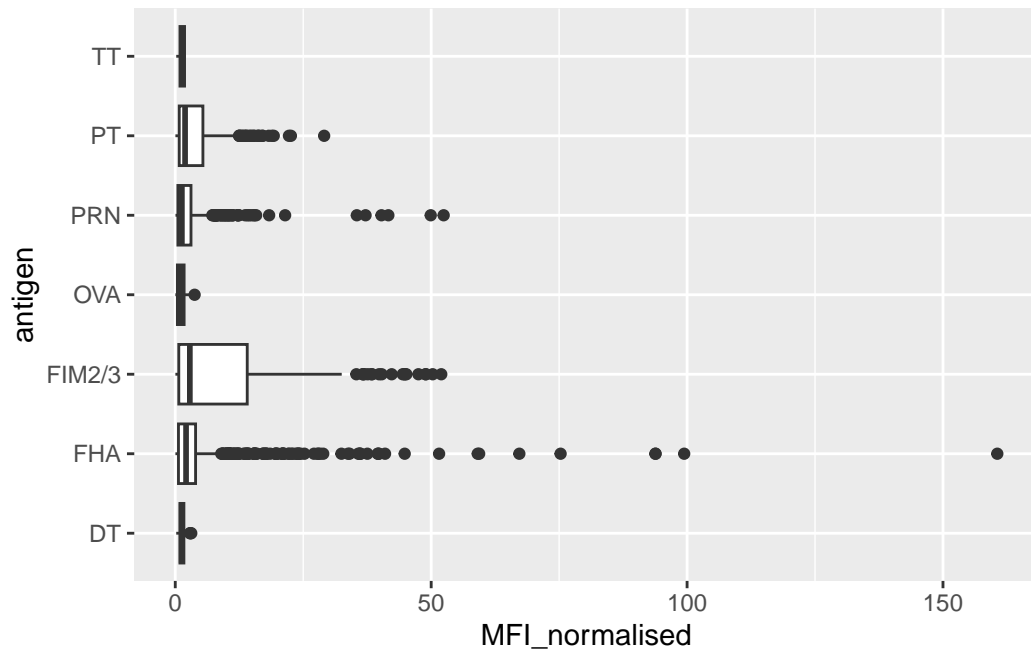
Q13. Complete the following code to make a summary boxplot of Ab titer levels (MFI) for all antigens:

```
ggplot(igg) +
  aes(MFI_normalised, antigen) +
  geom_boxplot()
```

```
ggplot(igg) +
  aes(MFI_normalised, antigen) +
  geom_boxplot() +
  xlim(0,75) +
  facet_wrap(vars(visit), nrow=2)
```

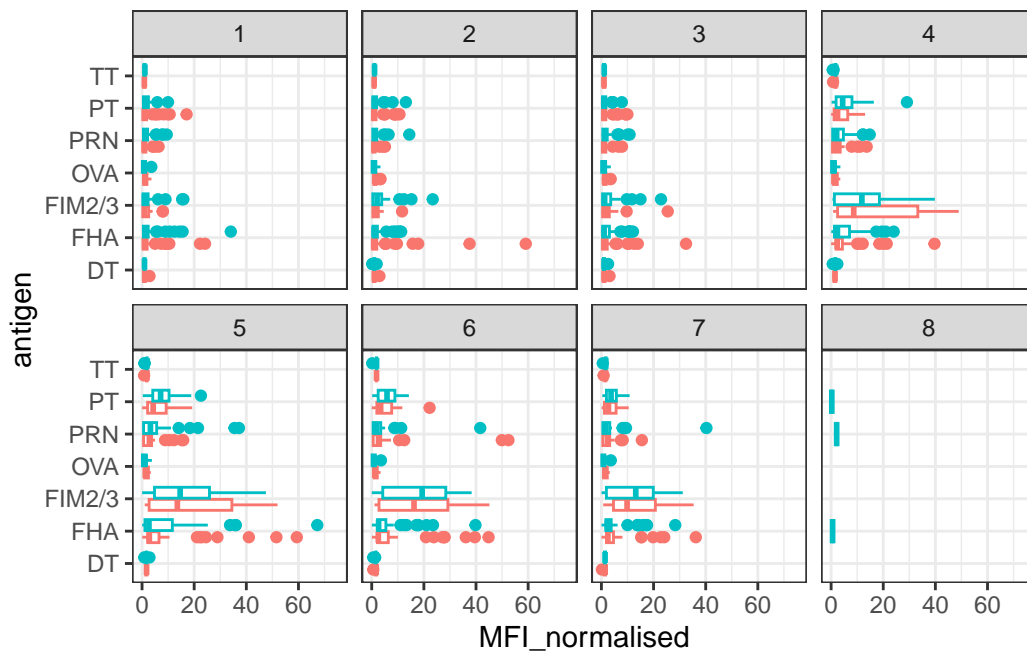Warning: Removed 5 rows containing non-finite values (`stat_boxplot()`).

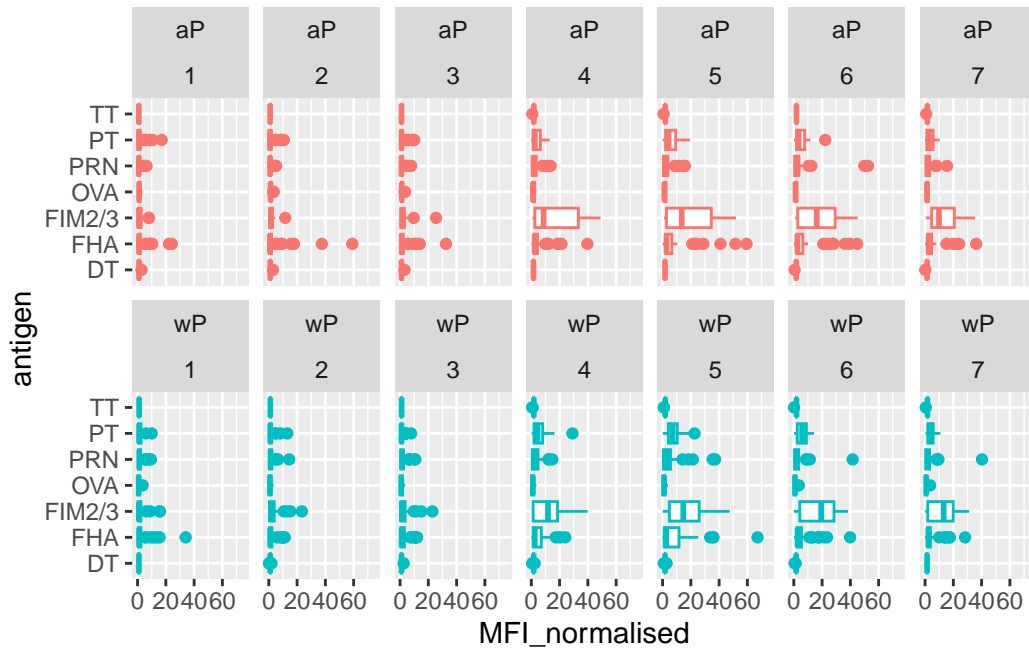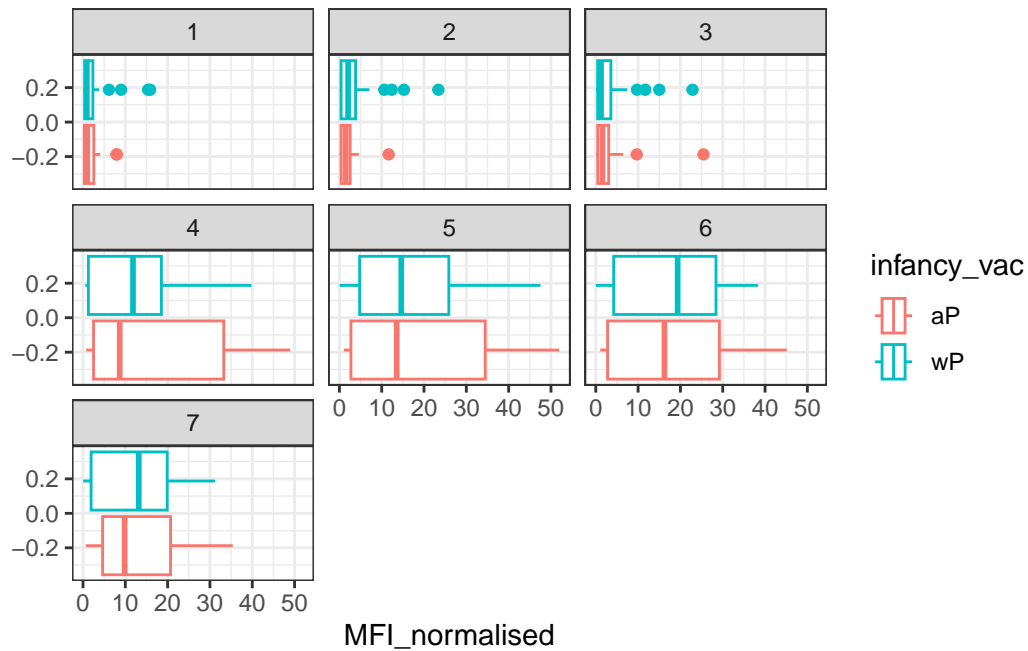Q14. What antigens show differences in the level of IgG antibody titers recognizing them over time? Why these and not others?

PT,PRN,FIM2/3,FHA show differences. These are surface molecules of pertussis that the patients are vaccinated against.

```
ggplot(igg) +
  aes(MFI_normalised, antigen, col=infancy_vac ) +
  geom_boxplot(show.legend = FALSE) +
  facet_wrap(vars(visit), nrow=2) +
  xlim(0,75) +
  theme_bw()
```

Warning: Removed 5 rows containing non-finite values (`stat_boxplot()`).

```
igg %>% filter(visit != 8) %>%
ggplot() +
  aes(MFI_normalised, antigen, col=infancy_vac ) +
  geom_boxplot(show.legend = FALSE) +
  xlim(0,75) +
  facet_wrap(vars(infancy_vac, visit), nrow=2)
```

Warning: Removed 5 rows containing non-finite values (`stat_boxplot()`).

Q15. Filter to pull out only two specific antigens for analysis and create a boxplot for each. You can chose any you like. Below I picked a "control" antigen ("OVA", that is not in our vaccines) and a clear antigen of interest ("PT", Pertussis Toxin, one of the key virulence factors produced by the bacterium B. pertussis).

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

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

MFI_normalised

Q16. What do you notice about these two antigens time courses and the PT data in particular?

PT levels clearly rise over time but not OVA. This trend is similar between wP and aP subjects.

Q17. Do you see any clear difference in aP vs. wP responses?

No

```
oops <- abdata %>% filter(antigen =="Fim2/3")
table(oops$dataset)
```

```
< table of extent 0 >
```

I want a time course of IgG MFI_normalised

```
#abdata$planned_day_relative_to_boost
```

```
igpt.21 <- abdata %>%filter(dataset == "2021_dataset",
                            isotype == "IgG",
                            antigen == "PT")
```

```
ggplot(igpt.21) +
  aes(planned_day_relative_to_boost,
      MFI_normalised,
      col=infancy_vac) +
  geom_point(alpha=0.6) +
  geom_line(aes(group=subject_id),linewidth=0.5,alpha=0.5) +
  geom_smooth(se= FALSE, span =0.4,linewidth=3)
```

`geom_smooth()` using method = 'loess' and formula = 'y ~ x'

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at -0.6

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 3.6

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: reciprocal condition number 1.8382e-16

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: There are other near singularities as well. 11364

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at -0.6

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 3.6

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
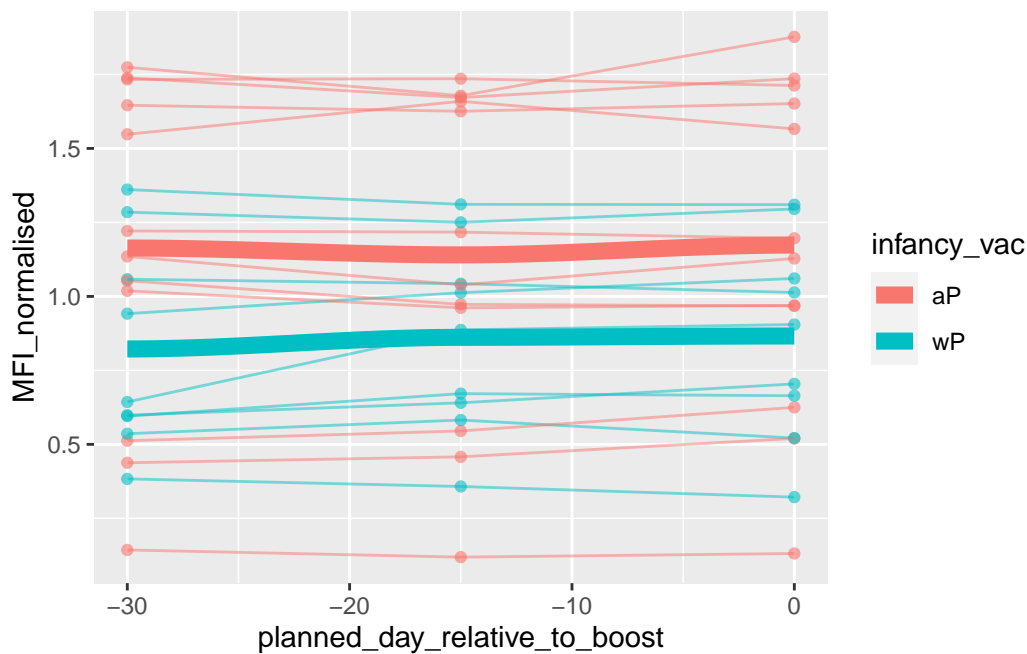: reciprocal condition number 1.4316e-16

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: There are other near singularities as well. 11364

```
igpt.22 <- abdata %>%filter(dataset == "2022_dataset",
                            isotype == "IgG",
                            antigen == "PT")
ggplot(igpt.22) +
  aes(planned_day_relative_to_boost,
      MFI_normalised,
      col=infancy_vac) +
  geom_point(alpha=0.6) +
  geom_line(aes(group=subject_id),linewidth=0.5,alpha=0.5) +
  geom_smooth(se= FALSE, span =0.4,linewidth=3)
```

`geom_smooth()` using method = 'loess' and formula = 'y ~ x'

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at -30.15

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 15.15

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: reciprocal condition number 0

```
Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: There are other near singularities as well. 229.52

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at -30.15

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 15.15

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: reciprocal condition number 0

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: There are other near singularities as well. 229.52
```



Q18. Does this trend look similar for the 2020 dataset?

```
igpt.20 <- abdata %>%filter(dataset == "2020_dataset",
                            isotype == "IgG",
                            antigen == "PT")
ggplot(igpt.20) +
```

```r
  aes(planned_day_relative_to_boost,
      MFI_normalised,
      col=infancy_vac) +
  geom_point(alpha=0.6) +
  geom_line(aes(group=subject_id),linewidth=0.5,alpha=0.5) +
  geom_smooth(se= FALSE, span =0.4,linewidth=3)
```

`geom_smooth()` using method = 'loess' and formula = 'y ~ x'

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at -0.6

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 3.6

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
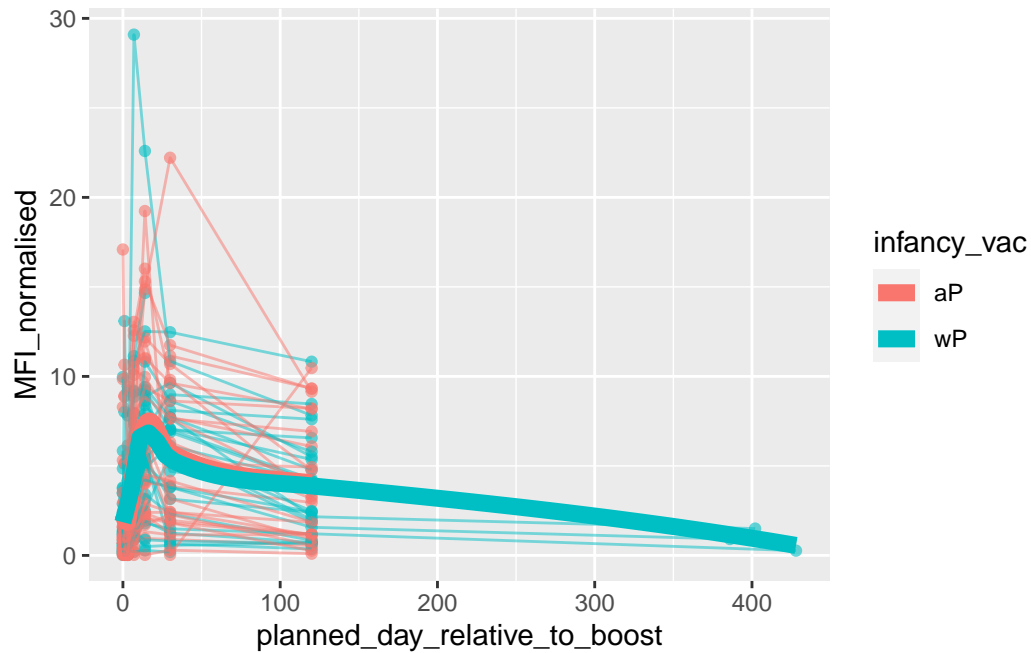: reciprocal condition number 2.9482e-16

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: There are other near singularities as well. 11364

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at -2.14

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 5.14

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: reciprocal condition number 4.7594e-16

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: There are other near singularities as well. 9
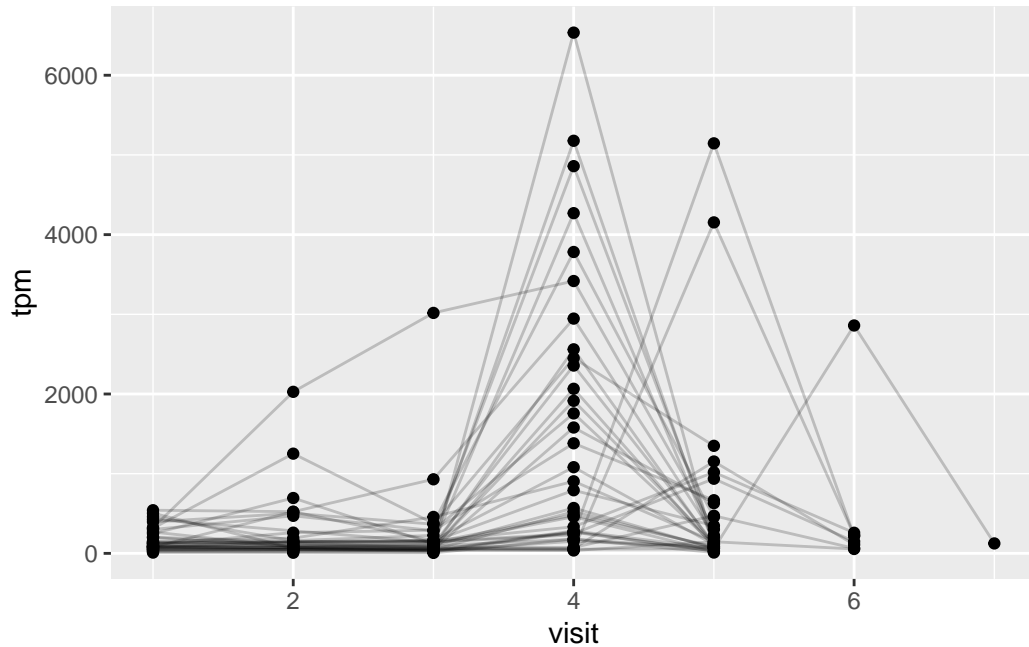
No.

# 5. Obtaining CMI-PB RNASeq data

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

```r
ssrna <- inner_join(rna, meta)
```

```
Joining with `by = join_by(specimen_id)`
```

Q19. Make a plot of the time course of gene expression for IGHG1 gene (i.e. a plot of visit vs. tpm).

```r
ggplot(ssrna) +
  aes(visit, tpm, group=subject_id) +
  geom_point() +
  geom_line(alpha=0.2)
```
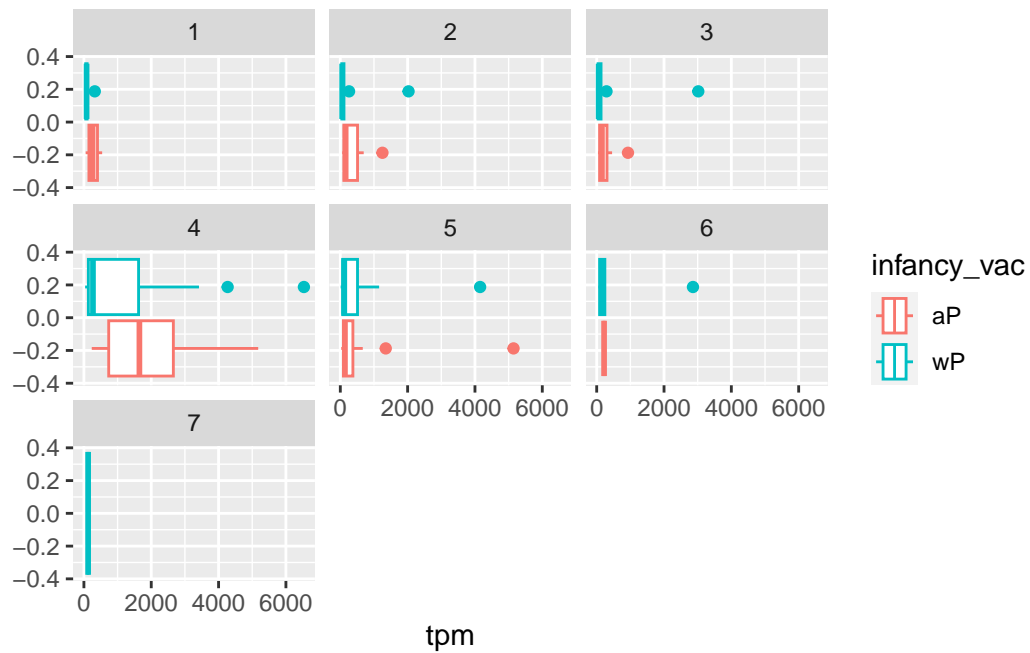
Q20.: What do you notice about the expression of this gene (i.e. when is it at it's maximum level)?

The maximum level is at the 4th visit.

Q21. Does this pattern in time match the trend of antibody titer data? If not, why not?

No.

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

```
ssrna %>%
  filter(visit==4) %>%
  ggplot() +
    aes(tpm, col=infancy_vac) + geom_density() +
    geom_rug()
```