# Class15:Pertussis and the CMI-PB project

Xiaoyan Wang(A16454055)

## 1. Investigating pertussis cases by year

**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.

```r
##install.packages("datapasta")
library(ggplot2)
```

```
Warning: package 'ggplot2' was built under R version 4.3.3
```

```r
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,2020L,2021L,2022L),
    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,
                                    32971,20762,17972,18975,15609,18617,6124,
                                    2116,3044)
)
```

```
head(cdc)
```

```
  Year No..Reported.Pertussis.Cases
1 1922                       107473
2 1923                       164191
3 1924                       165418
4 1925                       152003
5 1926                       202210
6 1927                       181411
```
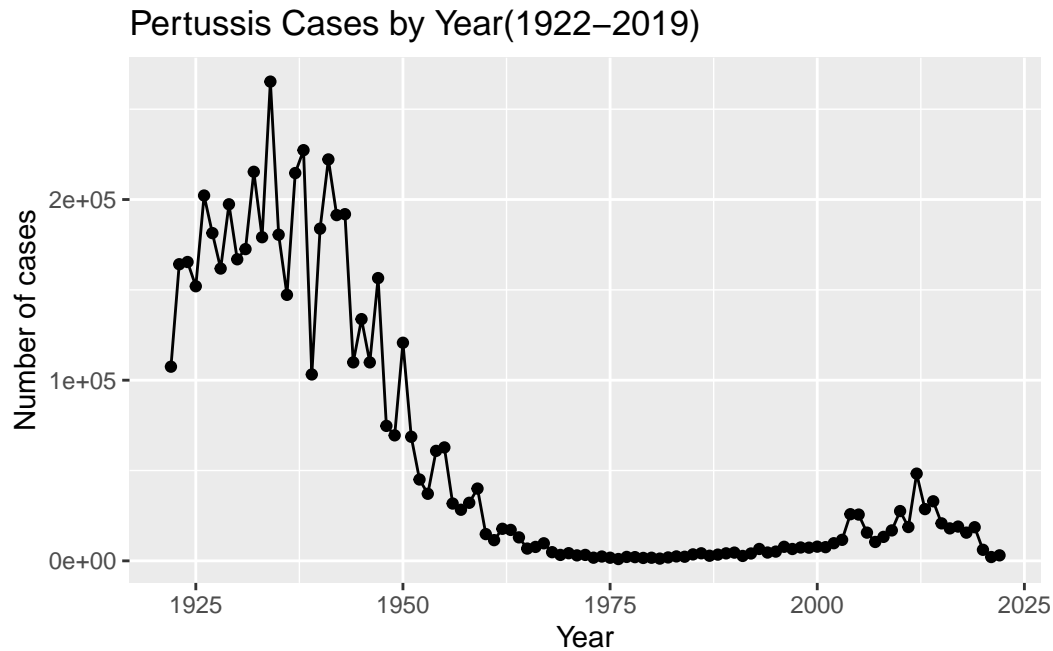
```
baseplot<- ggplot(cdc) +
  aes(x = Year, y = No..Reported.Pertussis.Cases) +
  geom_point() +
  geom_line() +
  labs(title = "Pertussis Cases by Year(1922-2019)", y = "Number of cases")
baseplot
```
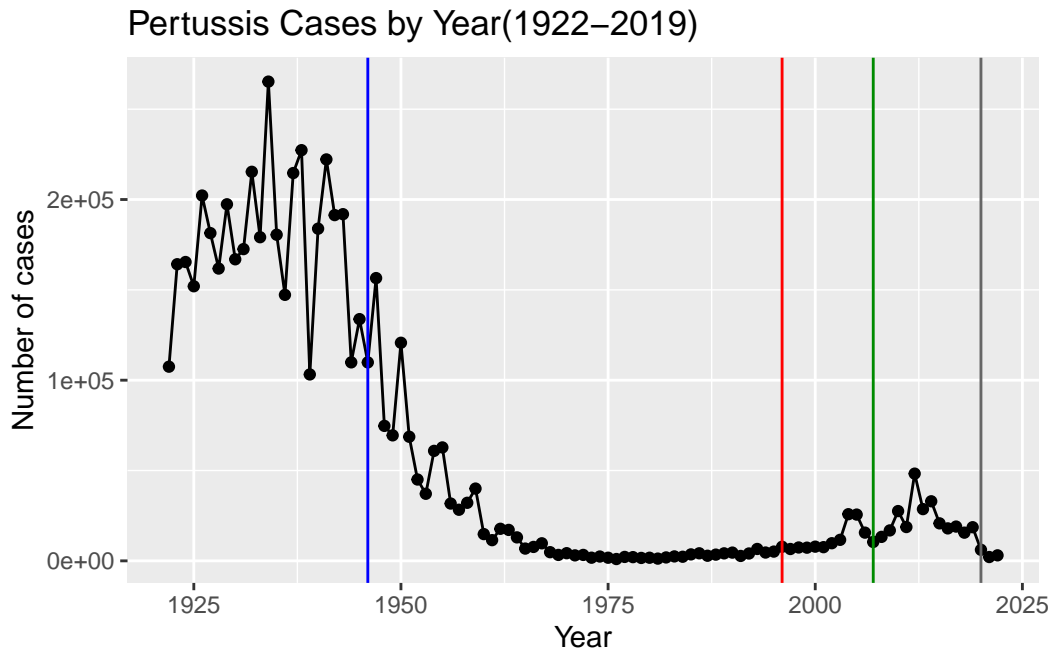
Pertussis Cases by Year(1922–2019)

## 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?

```
baseplot+
  geom_vline(xintercept = 1946, col = "blue") +
  geom_vline(xintercept = 1996, col = "red") +
  geom_vline(xintercept = 2020, col = "grey40") +
  geom_vline(xintercept = 2007, col = "green4")
```

## Pertussis Cases by Year(1922–2019)



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

The cases start to increase. No dramatic increase in the first few years but we start to see a big increase starting 2004, ~8-10 years after switching the vaccine.

There is a ~10 year log from aP roll out to increase cases numbers. This holds true of other countries like Japan, UK, etc.

**Key Question.** Why does the aP vaccine induced immunity wane faster than that of the wP vaccine?

## 3. Exploring CMI-PB data

```
# Allows us to read, write and process JSON data
library(jsonlite)
```

```
Warning: package 'jsonlite' was built under R version 4.3.3
```

```
subject <- read_json("https://www.cmi-pb.org/api/subject", simplifyVector = TRUE)
head(subject, 3)
```

4

```
  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
87 85
```

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

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

```
Female    Male
   112      60
```

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

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

```
        American Indian/Alaska Native Asian Black or African American
Female                             0    32                          2
Male                               1    12                          3

        More Than One Race Native Hawaiian or Other Pacific Islander
Female                  15                                         1
Male                    4                                          1

        Unknown or Not Reported White
Female                       14    48
Male                          7    32
```

Q: Does this do a good job of representing the US populus?

No it is not representative. There are some groups that are poorly representedd

```
library(lubridate)
```

Warning: package 'lubridate' was built under R version 4.3.3

Attaching package: 'lubridate'

The following objects are masked from 'package:base':

    date, intersect, setdiff, union

**Side-Note: Working with dates**

```
today()
```

[1] "2024-11-25"

```
today() - ymd("2000-01-01")
```

Time difference of 9095 days

```
time_length( today() - ymd("2000-01-01"),  "years")
```

[1] 24.90075

**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?

```
# Use todays date to calculate age in days
subject$age <- today() - ymd(subject$year_of_birth)
```

```
library(dplyr)
```

Warning: package 'dplyr' was built under R version 4.3.3

```
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
ap <- subject %>% filter(infancy_vac == "aP")

round( summary( time_length( ap$age, "years" ) ) )
```

```
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
     22      26      27      27      28      34
```

```r
# wP
wp <- subject %>% filter(infancy_vac == "wP")
round( summary( time_length( wp$age, "years" ) ) )
```

```
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
     22      32      34      36      39      57
```

**Q8.** Determine the age of all individuals at time of boost?

```r
int <- ymd(subject$date_of_boost) - ymd(subject$year_of_birth)
age_at_boost <- time_length(int, "year")
head(age_at_boost)
```

```
[1] 30.69678 51.07461 33.77413 28.65982 25.65914 28.77481
```

**Q9.** With the help of a faceted boxplot or histogram (see below), do you think these two groups are significantly different?

```
ggplot(subject) +
  aes(time_length(age, "year"),
      fill=as.factor(infancy_vac)) +
  geom_histogram(show.legend=FALSE) +
  facet_wrap(vars(infancy_vac), nrow=2) +
  xlab("Age in years")
```

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



```
# Complete the API URLs...
specimen <- read_json("https://www.cmi-pb.org/api/specimen", simplifyVector = TRUE)
titer <- read_json("https://www.cmi-pb.org/api/plasma_ab_titer", simplifyVector = TRUE)
head(titer)
```

```
  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
1 UG/ML                 2.096133
2 IU/ML                29.170000
3 IU/ML                 0.530000
4 IU/ML                 6.205949
5 IU/ML                 4.679535
6 IU/ML                 2.816431
```

**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:

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

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

```
dim(meta)
```

```
[1] 1503    14
```

```
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
1 14208 days
2 14208 days
3 14208 days
4 14208 days
5 14208 days
6 14208 days
```

**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] 52576      21
```

**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   5389 10117 10124 10124 10124
```

**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$visit)
```

```
   1     2     3     4     5     6     7     8     9    10    11    12
8280  8280  8420  6565  6565  6210  5810   815   735   686   105   105
```

## 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
1 14208 days
2 14208 days
3 14208 days
4 15304 days
5 15304 days
6 15304 days
```
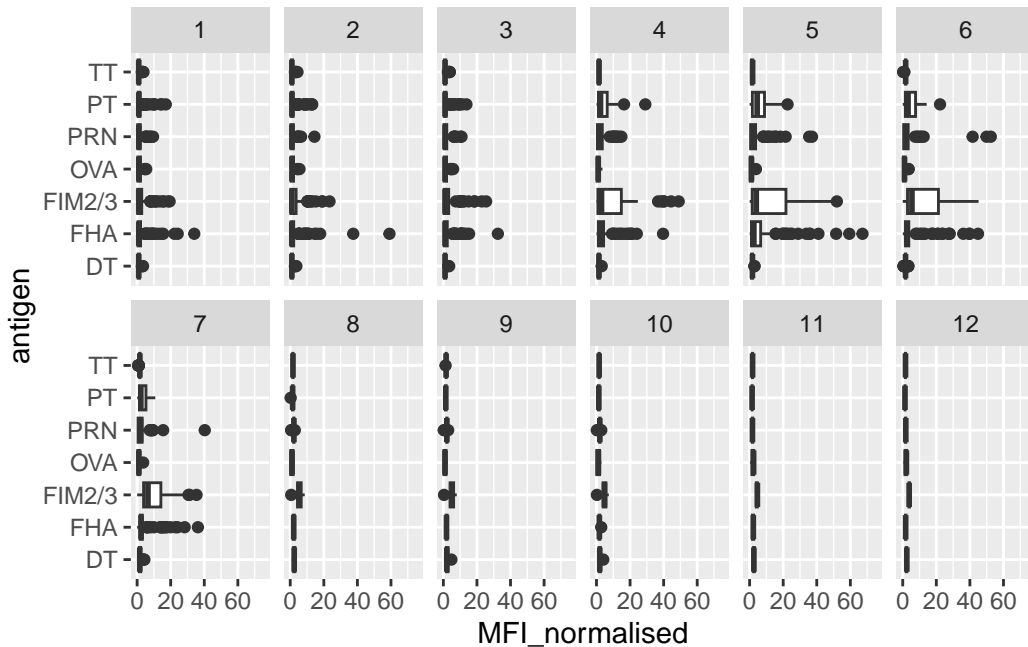
**Q13.** Complete the following code to make a summary boxplot of Ab titer levels (MFI) for all antigens:

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

Warning: Removed 5 rows containing non-finite outside the scale range
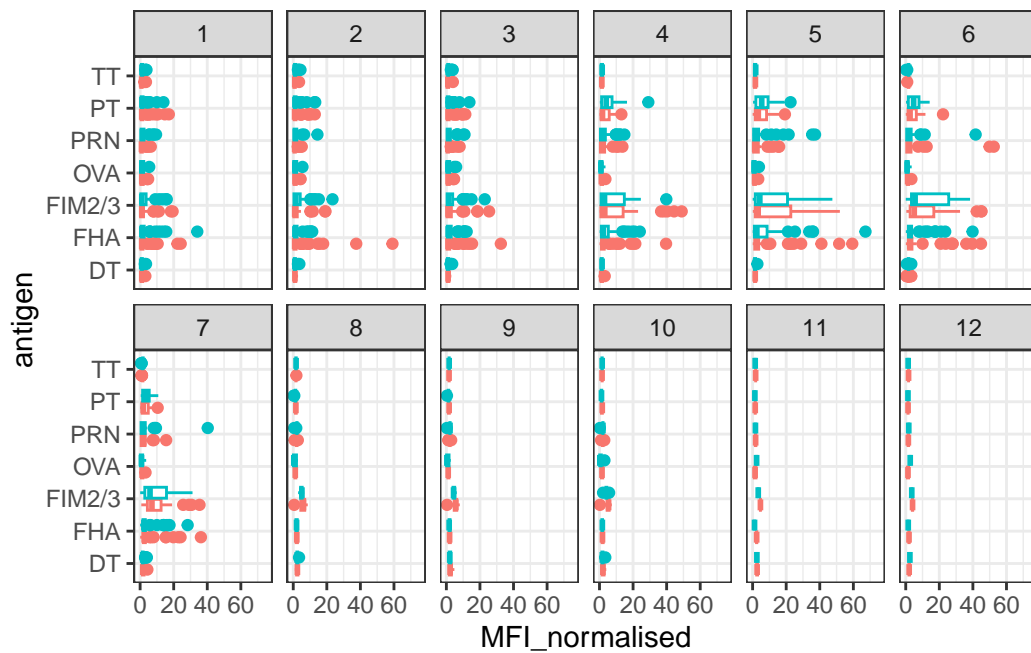(`stat_boxplot()`).



**Q14.** What antigens show differences in the level of IgG antibody titers recognizing them over time? Why these and not others?
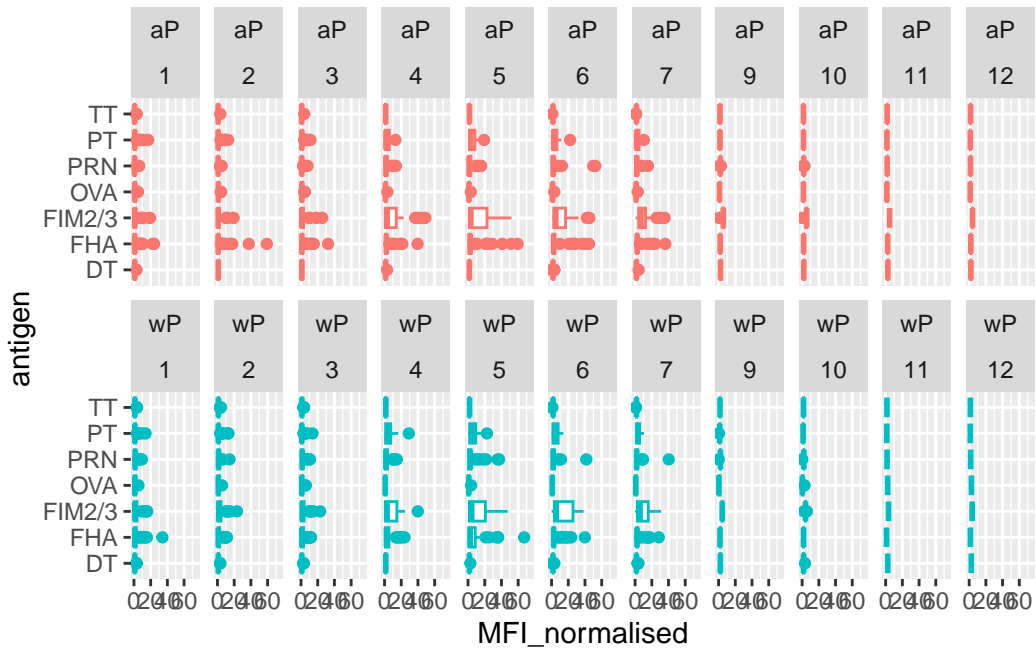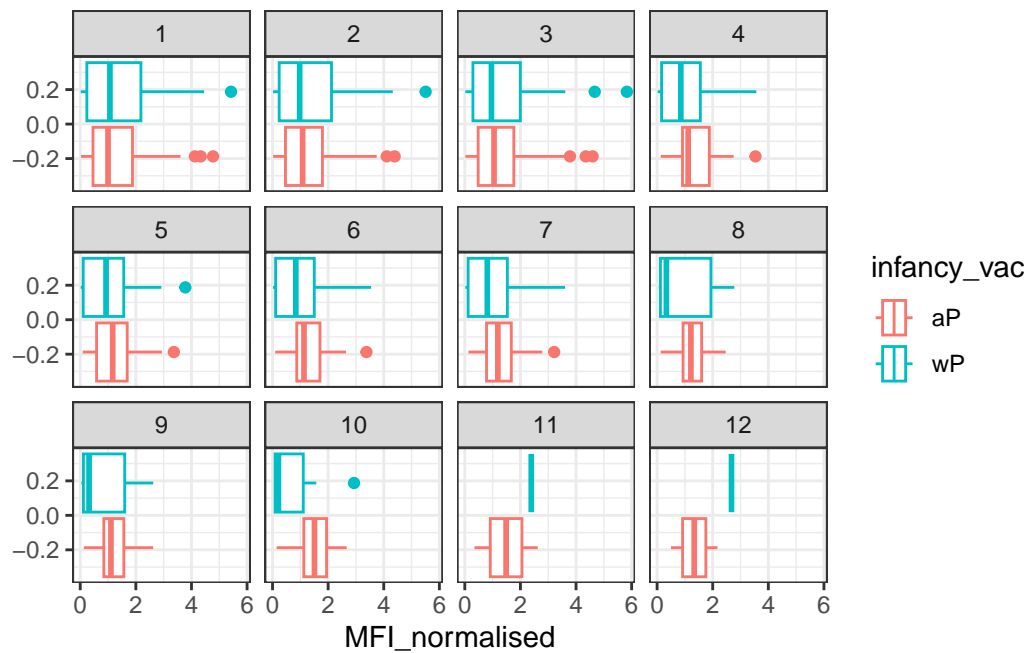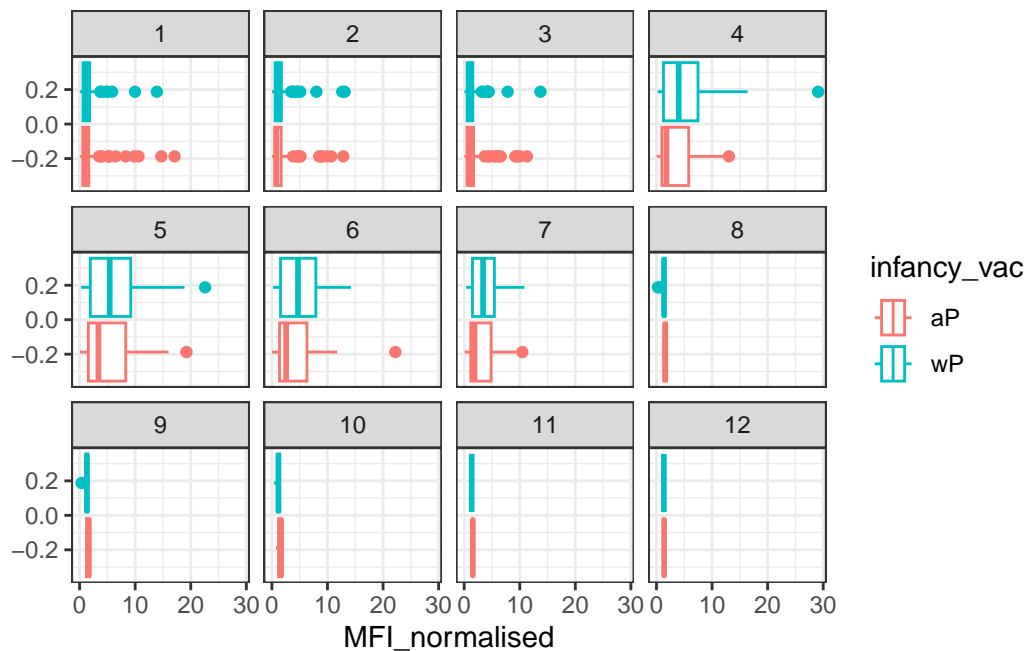
PT, PRN, FIM2/3 and FHA

```
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 outside the scale range
(`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 outside the scale range
(`stat_boxplot()`).

13

**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=="PT") %>%
  ggplot() +
  aes(MFI_normalised, col=infancy_vac) +
  geom_boxplot(show.legend = TRUE) +
  facet_wrap(vars(visit)) +
  theme_bw()
```

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

PT seems to have a rise since day 4 and drops back to a very low level on week 8.

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

Yes in OVA but not in PT. The wP in OVA data seems to have a smaller median compare to the aP(if consider <day10)

```
abdata.21 <- abdata %>% filter(dataset == "2021_dataset")

abdata.21 %>%
  filter(isotype == "IgG",  antigen == "PT") %>%
  ggplot() +
    aes(x=planned_day_relative_to_boost,
        y=MFI_normalised,
        col=infancy_vac,
        group=subject_id) +
    geom_point() +
    geom_line() +
    geom_vline(xintercept=0, linetype="dashed") +
    geom_vline(xintercept=14, linetype="dashed") +
  labs(title="2021 dataset IgG PT",
       subtitle = "Dashed lines indicate day 0 (pre-boost) and 14 (apparent peak levels)")
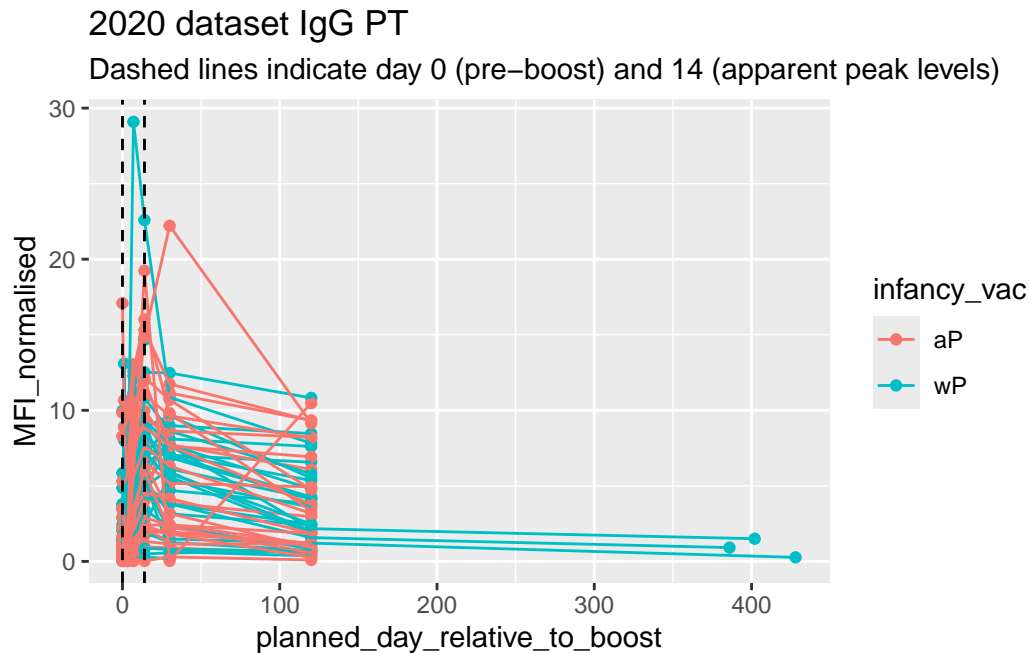```

# 2021 dataset IgG PT
## Dashed lines indicate day 0 (pre−boost) and 14 (apparent peak levels)



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

```r
abdata.20 <- abdata %>% filter(dataset == "2020_dataset")

abdata.20 %>%
  filter(isotype == "IgG",  antigen == "PT") %>%
  ggplot() +
    aes(x=planned_day_relative_to_boost,
        y=MFI_normalised,
        col=infancy_vac,
        group=subject_id) +
    geom_point() +
    geom_line() +
    geom_vline(xintercept=0, linetype="dashed") +
    geom_vline(xintercept=14, linetype="dashed") +
  labs(title="2020 dataset IgG PT",
        subtitle = "Dashed lines indicate day 0 (pre-boost) and 14 (apparent peak levels)")
```

## 2020 dataset IgG PT
### Dashed lines indicate day 0 (pre−boost) and 14 (apparent peak levels)



Not really. The 2020 dataset has an earlier peak of wP and a later peak of aP appearing, but for 2021 the aP and wP both have a peak at 14 days. Plus, the increase of MFI is more dramatic in 2020.

## 5. Obtaining CMI-PB RNASeq data

```
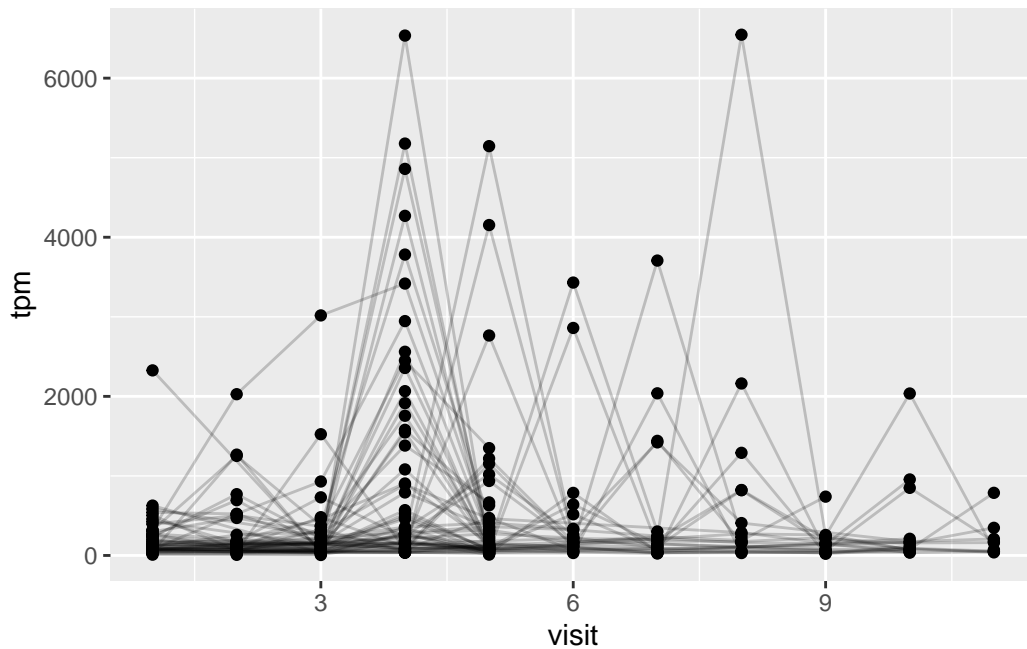url <- "https://www.cmi-pb.org/api/v2/rnaseq?versioned_ensembl_gene_id=eq.ENSG00000211896.7"

rna <- read_json(url, simplifyVector = TRUE)
```

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

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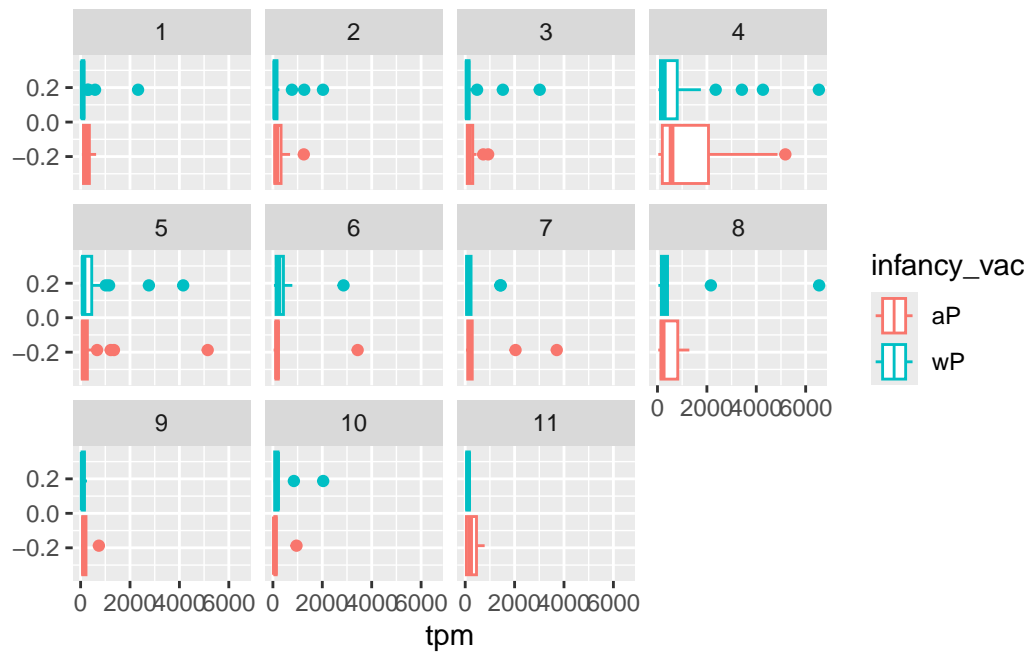  geom_point() +
  geom_line(alpha=0.2)
```



**Q20.**: What do you notice about the expression of this gene (i.e. when is it at it's maximum level)?

The expression of tpm raises to peak at visit 4, except a different one that raises to the peak twice, one in visit 4 and the other in visit 8. The general trend is to decrease after visit

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

This aligns with the antibody titer data in q15. As the maximum antibody was produced in ~week 5, the gene expression have to increase prior to that to encode antibody proteins. Therefore, it is reasonable that the gene expression was first raised, and then leading to a increase in antibodies.

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

19

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