

Pertussis Mini Project

Nhi To (PID: A18053310)

Pertussis (a.k.a Whooping Cough) is a deadly lung infection caused by the bacteria B. Pertussis

THE CDC tracks Pertussis cases around the US. <https://tinyurl.com/pertussiscdc>

We can “scrape” this data using the R **datapasta** package

```
head(cdc)
```

```
  year  cases
1 1922 107473
2 1923 164191
3 1924 165418
4 1925 152003
5 1926 202210
6 1927 181411
```

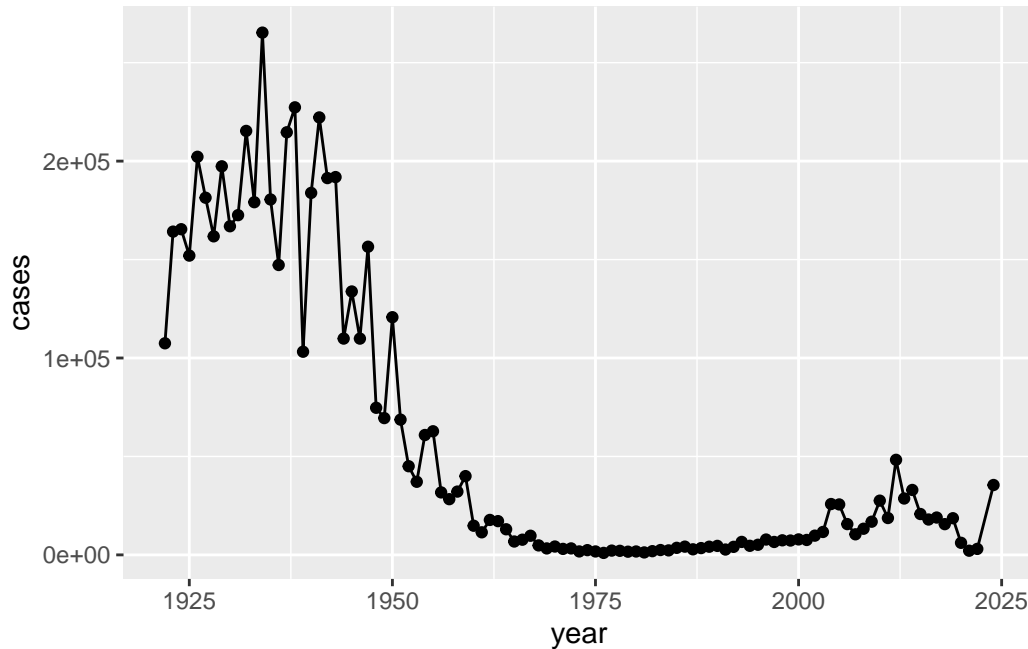
Q1: With the help of the R “addin” package datapasta assign the CDC pertussis case number data to a data frame called cdc and use ggplot to make a plot of cases numbers over time.

ANSWER:

```
library(ggplot2)

#data
#aes
#geoms

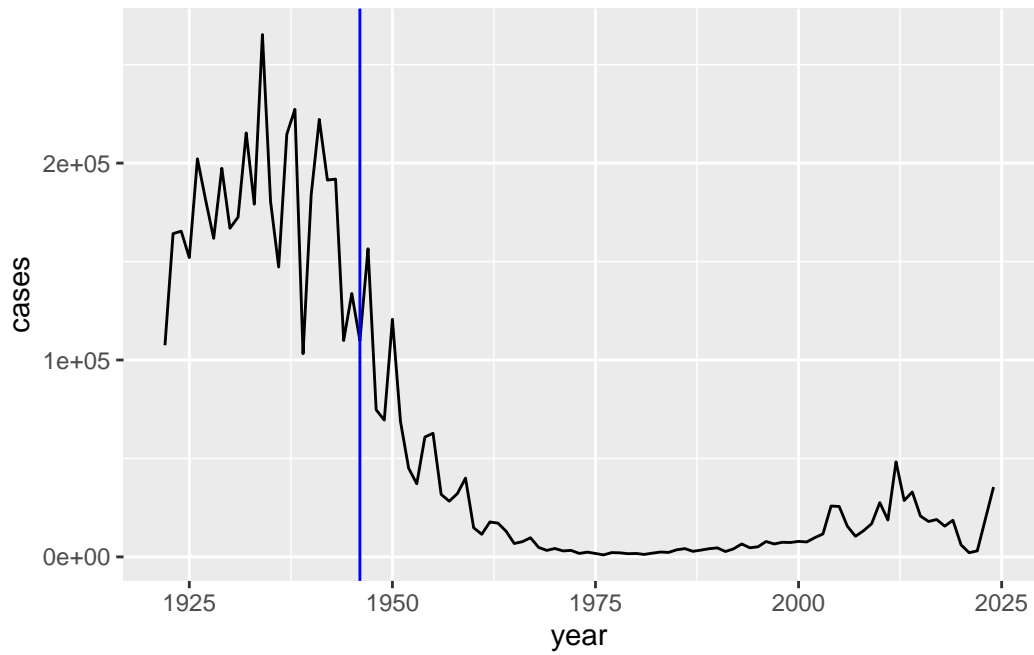
ggplot(cdc) +
  aes(year, cases) +
  geom_point() +
  geom_line() +
  labs(cdc)
```



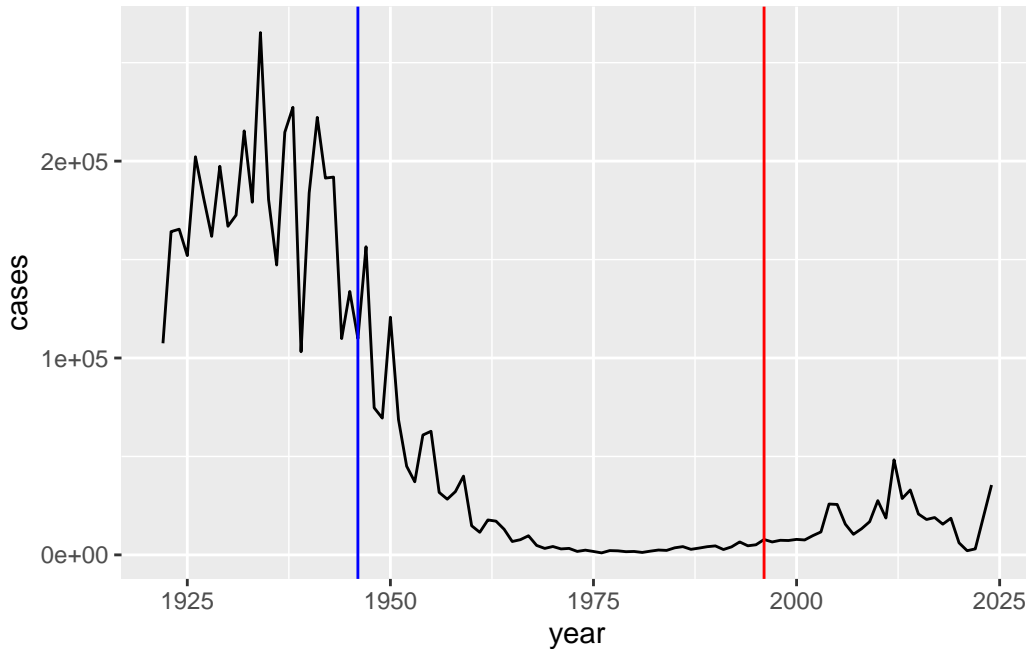
Q2: Using the ggplot `geom_vline()` function add lines to your previous plot for the 1946 introduction of the wP vaccine and the 1996 switch to aP vaccine (see example in the hint below). What do you notice?

ANSWER: In regards to the introduction of the wP vaccine, there is a tremendous drop in cases, and this drop is maintained until 1996. I noticed that there is a delay in vaccination benefits after the switch to aP vaccine compared to the wP vaccine as we noticed the rise in cases. In addition, there is also concern that aP vaccine benefits are waning, and is still be studied.

```
ggplot(cdc) +
  aes(year, cases) +
  geom_line() +
  geom_vline(xintercept= 1946, col= "blue")
```



```
ggplot(cdc) +  
  aes(year, cases) +  
  geom_line() +  
  geom_vline(xintercept= 1946, col= "blue") +  
  geom_vline(xintercept=1996, col="red")
```

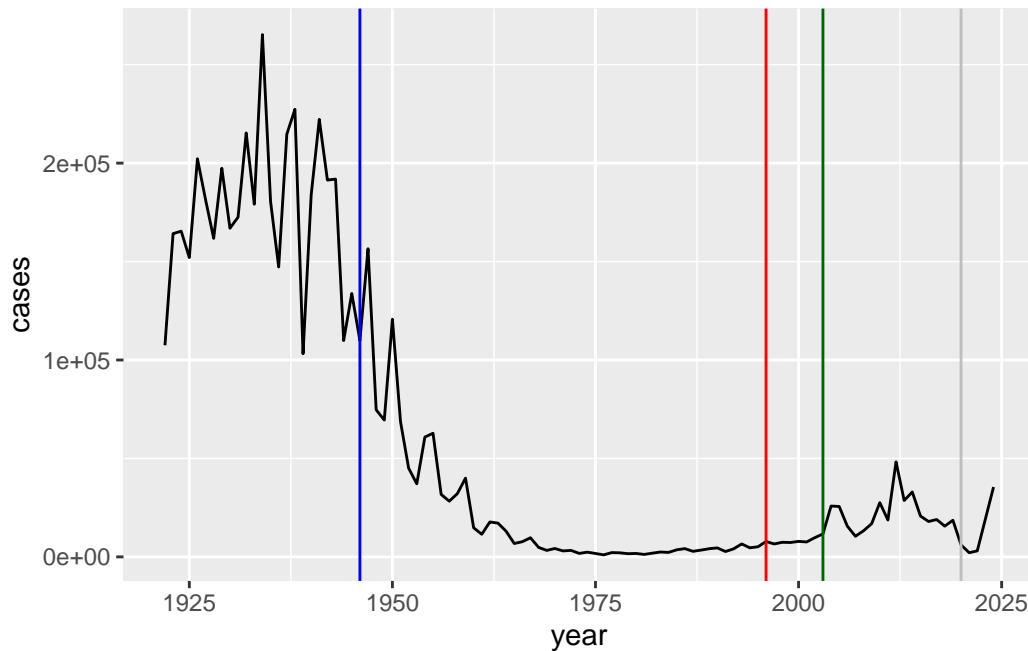


Q3. Describe what happened after the introduction of the aP vaccine? Do you have a possible explanation for the observed trend?

ANSWER: After the introduction of the aP vaccine, the pertussis cases began increasing again, as we seen in the graph above. aP vaccine was introduced in 1996, and starting around the 2000, the line shows an increase in cases. This can possibly be due to PCR-based testing being better, and now being able to detect more cases. Another reasoning could be the increasing doubt of vaccine benefits due to anti-vaccine propoganda shows form the past. Another reasoning could be that pertusssis has evolved over time, and has thus developed some amount of resistancy to the vaccine.

Adding a line for COVID

```
ggplot(cdc) +
  aes(year, cases) +
  geom_line() +
  geom_vline(xintercept= 1946, col= "blue") +
  geom_vline(xintercept=1996, col="red") +
  geom_vline(xintercept=2020, col="gray") +
  geom_vline(xintercept=2003, col="darkgreen")
```



There were high case numbers for the first wP (whole-cell) vaccine roll out in 1946 then a rapid decline in case numbers until 2004 when we have our first large-scale outbreaks of pertussis again. There is also a notable COVID related dip and recent rapid rise

Q. What is different about the immune response to infection if you had an older wP vaccine vs the newer aP vaccine?

Computational Models of Immunity Pertussis Boost (CMI-PB)

The CMI-PB project aims to address this key question. What is the difference between aP and wP individuals.

We can get all data from this ongoing project via JSON API calls

For this we will use the **jsonlite** package

```
library(jsonlite)
subject <- read_json("https://www.cmi-pb.org/api/v5_1/subject",
                     simplifyVector=TRUE)

head(subject)
```

```
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
4	4	wP	Male Not Hispanic or Latino Asian
5	5	wP	Male Not Hispanic or Latino Asian
6	6	wP	Female Not Hispanic or Latino 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
4	1988-01-01	2016-08-29	2020_dataset
5	1991-01-01	2016-08-29	2020_dataset
6	1988-01-01	2016-10-10	2020_dataset

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

ANSWER: 87 aP and 85 wP

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

```
aP wP
87 85
```

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

ANSWER: 112 females and 60 males

```
nrow(subject)
```

```
[1] 172
```

```
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...)?

ANSWER:

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

	Female	Male
American Indian/Alaska Native	0	1
Asian	32	12
Black or African American	2	3
More Than One Race	15	4
Native Hawaiian or Other Pacific Islander	1	1
Unknown or Not Reported	14	7
White	48	32

Side-note: working with dates

```
library(lubridate)
```

Attaching package: 'lubridate'

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

date, intersect, setdiff, union

```
today()
```

```
[1] "2025-03-09"
```

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

Time difference of 9199 days

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

```
[1] 25.18549
```

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?

ANSWER: The average age of wP individuals is 36 while the average age of aP individuals is 27. the difference is around 10 years, which is pretty significant.

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

```
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

```
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

```
# 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?

ANSWER:

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

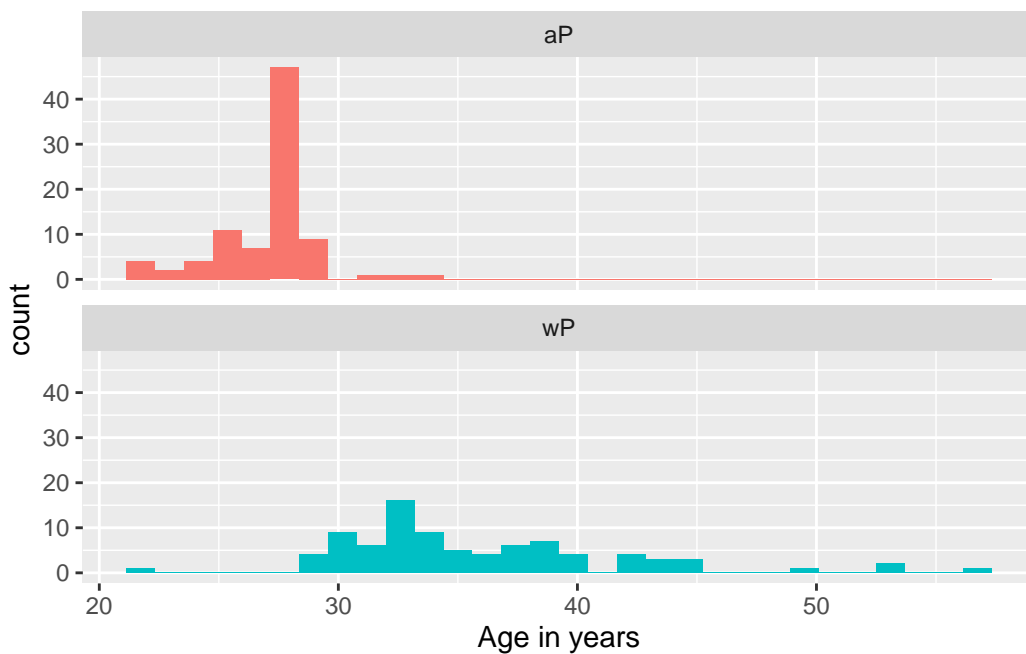
```
[1] 30.69678 51.07461 33.77413 28.65982 25.65914 28.77481
```


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

ANSWER: These two graphs highlight that the two groups are significantly different. In the aP graph, it is heavily skewed right with a peak around age 25 (years). The graph for wP shows that there is not really a peak and that there is a more even distribution. In addition, it is observed wP has much more older ages getting vaccinated compared to aP.

```
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`.



```
# Or use wilcox.test()  
x <- t.test(time_length( wp$age, "years" ),  
            time_length( ap$age, "years" ))  
  
x$p.value
```

[1] 2.372101e-23

Obtain more data from CMI-PB

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:

ANSWER:

```
library(jsonlite)
specimen <- read_json("https://www.cmi-pb.org/api/v5_1/specimen",
                      simplifyVector=TRUE)
ab_data <- read_json("https://www.cmi-pb.org/api/v5_1/plasma_ab_titer",
                    simplifyVector=TRUE)
```

```
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

```
head(ab_data)
```

	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	

I now have 3 tables of data from CMI-PB: 'subject', 'specimen', and 'ab_data'. I need to "join" these tables so I will have all the info I need to work with.

For this, we will use the 'inner_join()' function from the **dplyr** package.

```
library(dplyr)

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 14312 days
2 14312 days
3 14312 days
4 14312 days
5 14312 days
6 14312 days

```

```
dim(subject)
```

```
[1] 172  9
```

```
dim(specimen)
```

```
[1] 1503  6
```

```
dim(meta)
```

```
[1] 1503 14
```

Now we can join our 'ab_data' table to 'meta' so we have all the info we need about antibody levels.

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.

ANSWER:

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

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

```
head(abdata)
```

```

specimen_id subject_id actual_day_relative_to_boost
1           1           1                        -3
2           1           1                        -3
3           1           1                        -3
4           1           1                        -3
5           1           1                        -3
6           1           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 isotype is_antigen_specific antigen MFI MFI_normalised
1 14312 days IgE FALSE Total 1110.21154 2.493425
2 14312 days IgE FALSE Total 2708.91616 2.493425
3 14312 days IgG TRUE PT 68.56614 3.736992
4 14312 days IgG TRUE PRN 332.12718 2.602350
5 14312 days IgG TRUE FHA 1887.12263 34.050956
6 14312 days 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

```

```
dim(abdata)
```

```
[1] 61956 21
```

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

ANSWER:

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

```
  IgE   IgG  IgG1  IgG2  IgG3  IgG4
6698  7265 11993 12000 12000 12000
```

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

ANSWER: I noticed that the the number of rows for the most “recent” (2023) dataset has increased by double since the last dataset in 2022

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

```
2020_dataset 2021_dataset 2022_dataset 2023_dataset
          31520           8085           7301           15050
```

Q. How many different antibody isotypes are there in this dataset?

```
length(abdata$isotype)
```

```
[1] 61956
```

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

```
61956
  1
```

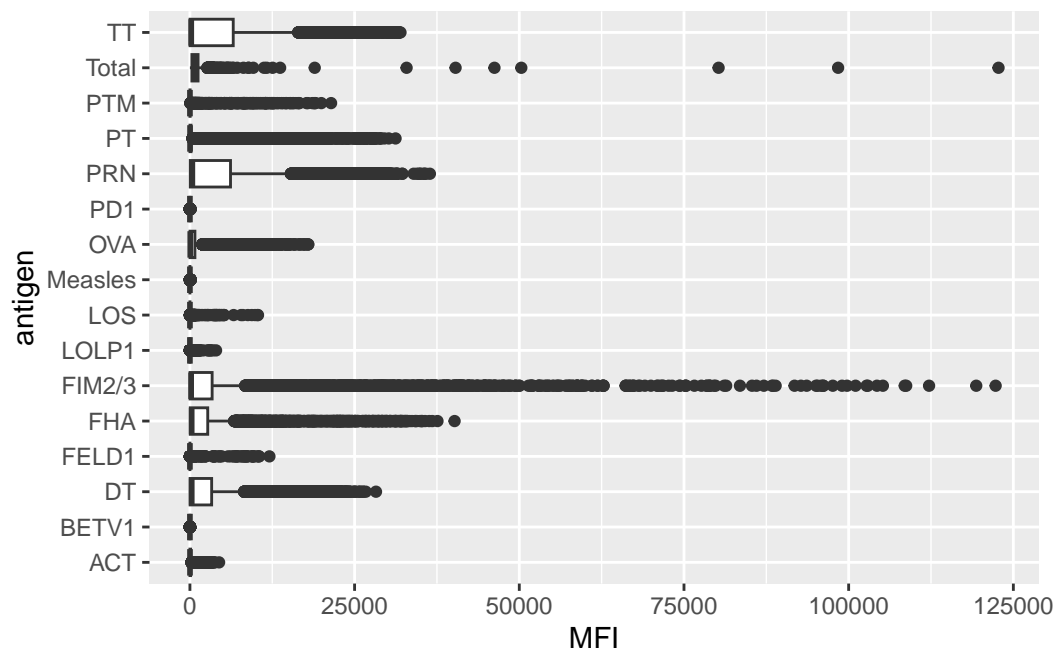
```
table(abdata$antigen)
```

```
  ACT  BETV1    DT  FELD1    FHA  FIM2/3  LOLP1    LOS Measles    OVA
1970  1970  6318  1970  6712  6318  1970  1970  1970  6318
 PD1   PRN    PT   PTM  Total    TT
1970  6712  6712  1970   788  6318
```

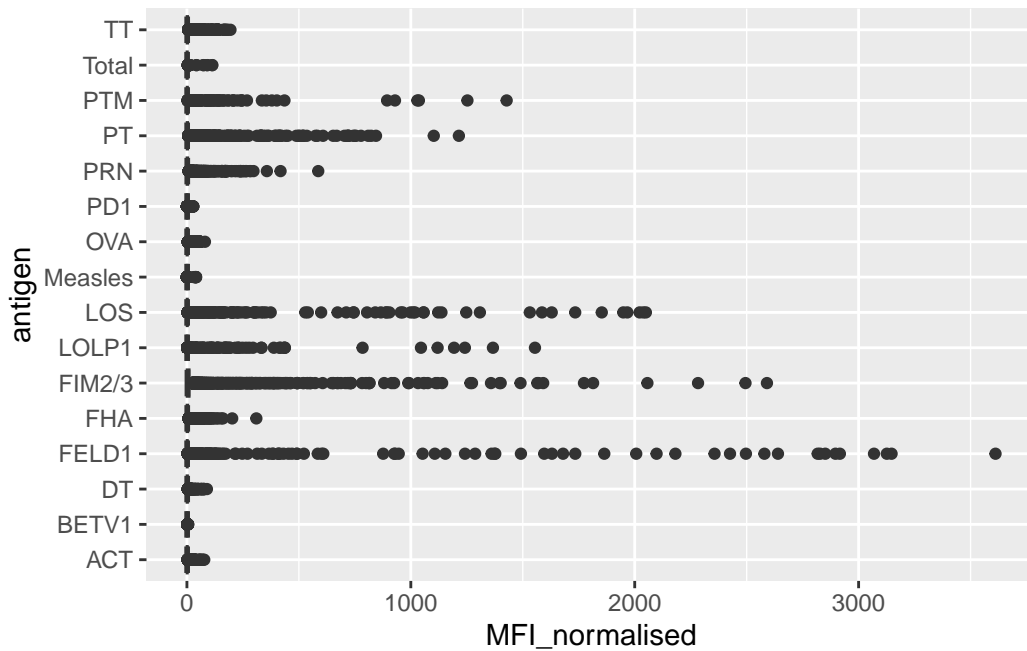
I want a plot of antigen levels across the whole dataset

```
ggplot(abdata)+  
  aes(MFI, antigen) +  
  geom_boxplot()
```

Warning: Removed 1 row containing non-finite outside the scale range
(`stat_boxplot()`).



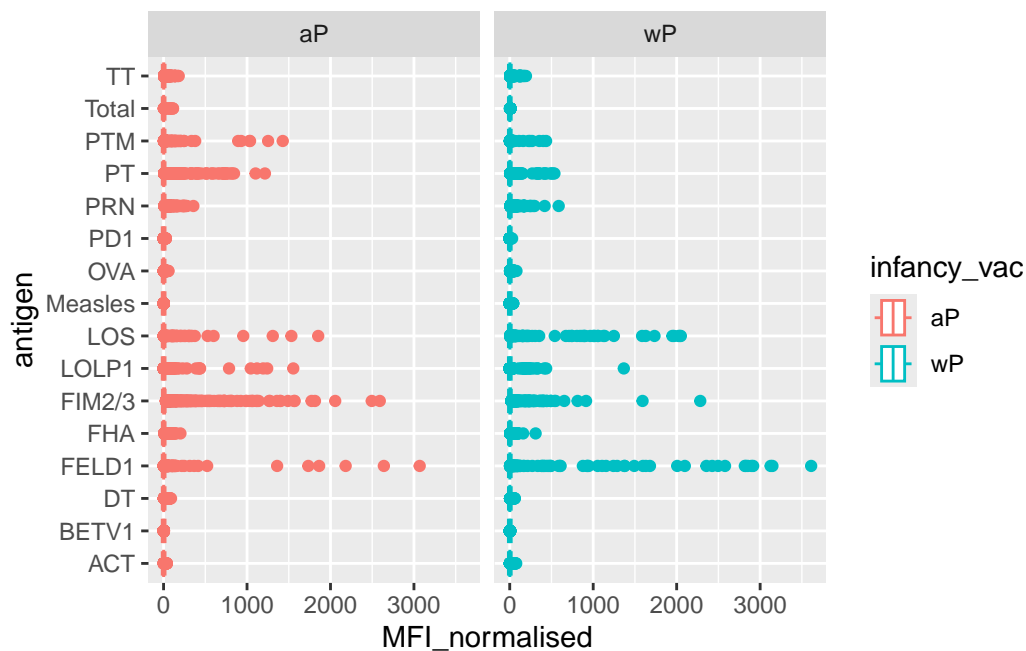
```
ggplot(abdata)+  
  aes(MFI_normalised, antigen) +  
  geom_boxplot()
```



Antigens like FIM2/3, PT, FELD1 have quite a large range of values Others like Measles don't show much activity

Q. Are there differences at this whole-dataset level because aP and wP?

```
ggplot(abdata)+
  aes(MFI_normalised, antigen, col=infancy_vac) +
  geom_boxplot() +
  facet_wrap(~infancy_vac)
```

4. Examine IgG Ab titer levels

For this I need to select out just isotype IgG

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

	specimen_id	subject_id	actual_day_relative_to_boost			
1	1	1	-3			
2	1	1	-3			
3	1	1	-3			
4	2	1	1			
5	2	1	1			
6	2	1	1			
	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	1	Blood	2	wP	Female	
5	1	Blood	2	wP	Female	
6	1	Blood	2	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	isotype	is_antigen_specific	antigen	MFI	MFI_normalised
1	14312 days	IgG	TRUE	PT	68.56614	3.736992
2	14312 days	IgG	TRUE	PRN	332.12718	2.602350
3	14312 days	IgG	TRUE	FHA	1887.12263	34.050956
4	14312 days	IgG	TRUE	PT	41.38442	2.255534
5	14312 days	IgG	TRUE	PRN	174.89761	1.370393
6	14312 days	IgG	TRUE	FHA	246.00957	4.438960

	unit	lower_limit_of_detection
1	IU/ML	0.530000
2	IU/ML	6.205949
3	IU/ML	4.679535
4	IU/ML	0.530000
5	IU/ML	6.205949
6	IU/ML	4.679535

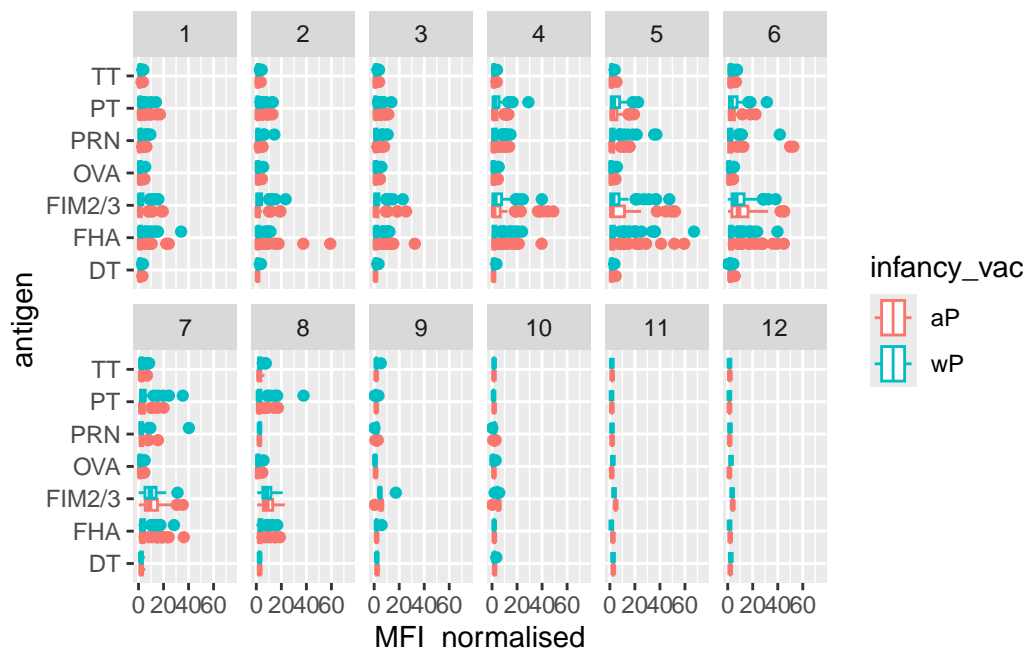
Q13. Complete the following code to make a summary boxplot of Ab titer levels (MFI) for all antigens:

ANSWER:

A overview boxplot:

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

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

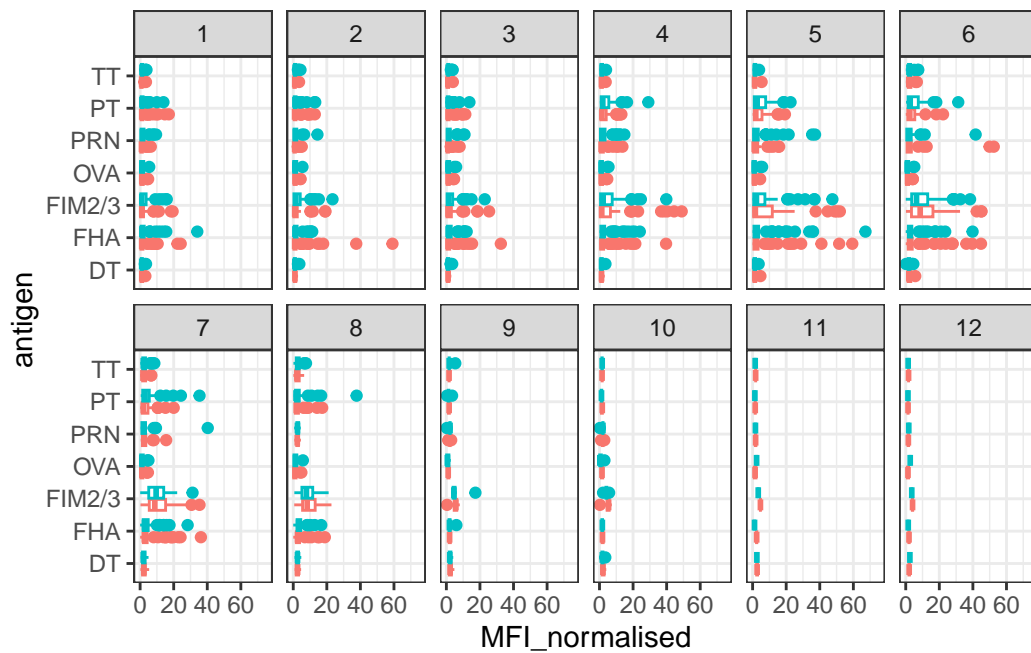


Q14. What antigens show differences in the level of IgG antibody titers recognizing them over time? Why these and not others?

ANSWER: PT and FIM2/3 antigens show the most noticeable differences because their levels are significantly higher than the other antigens in some visits, such as them peaking at visit #5. There should be differences in the levels of IgG antibody titers because antibody responses will change once the vaccination is injected. The control antigens are not expected to have noticeable differences, because they are controls that aren't incorporated into the vaccines.

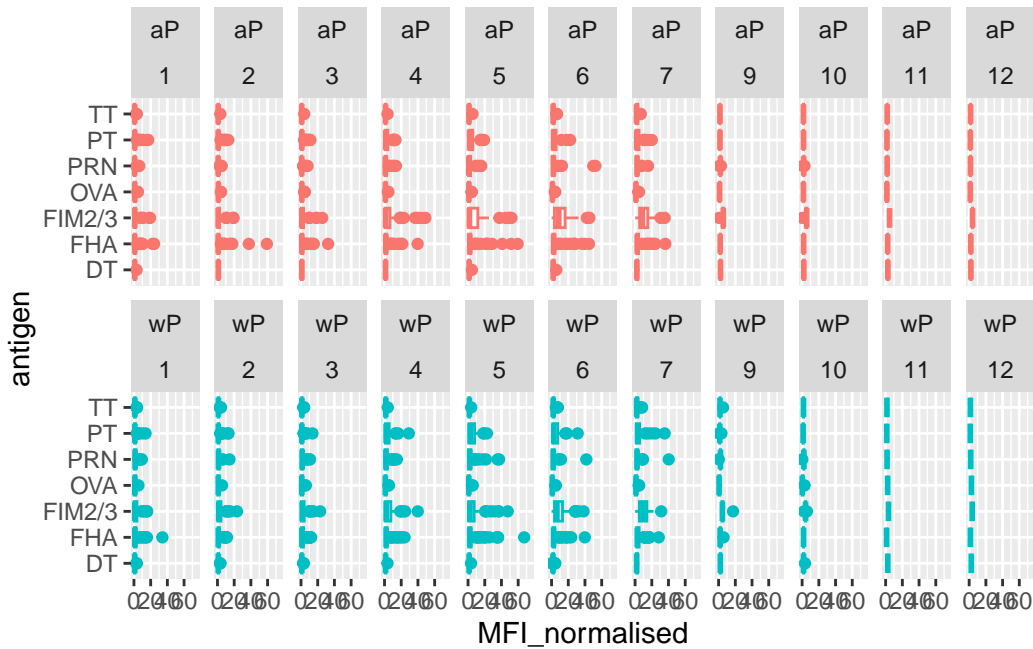
```
ggplot(igg) +
  aes(MFI_normalised, antigen, col=infancy_vac ) +
  geom_boxplot(show.legend = FALSE) +
  facet_wrap(vars(visit), nrow=2) +
  xlim(0,75) +
  theme_bw()
```

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)
```

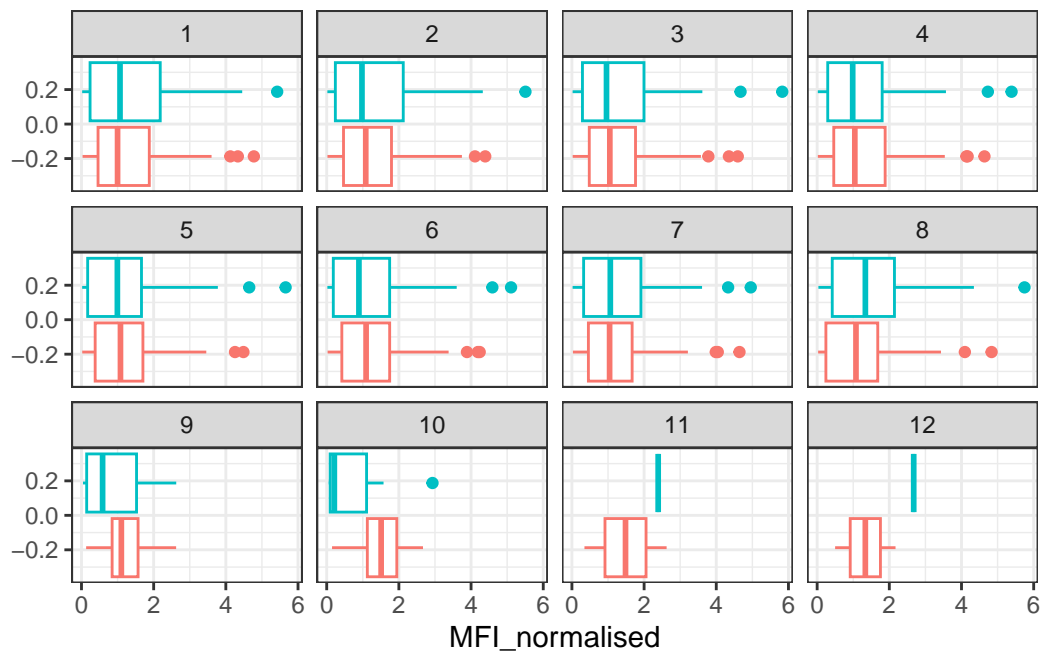
Warning: Removed 5 rows containing non-finite outside the scale range (`stat_boxplot()`).



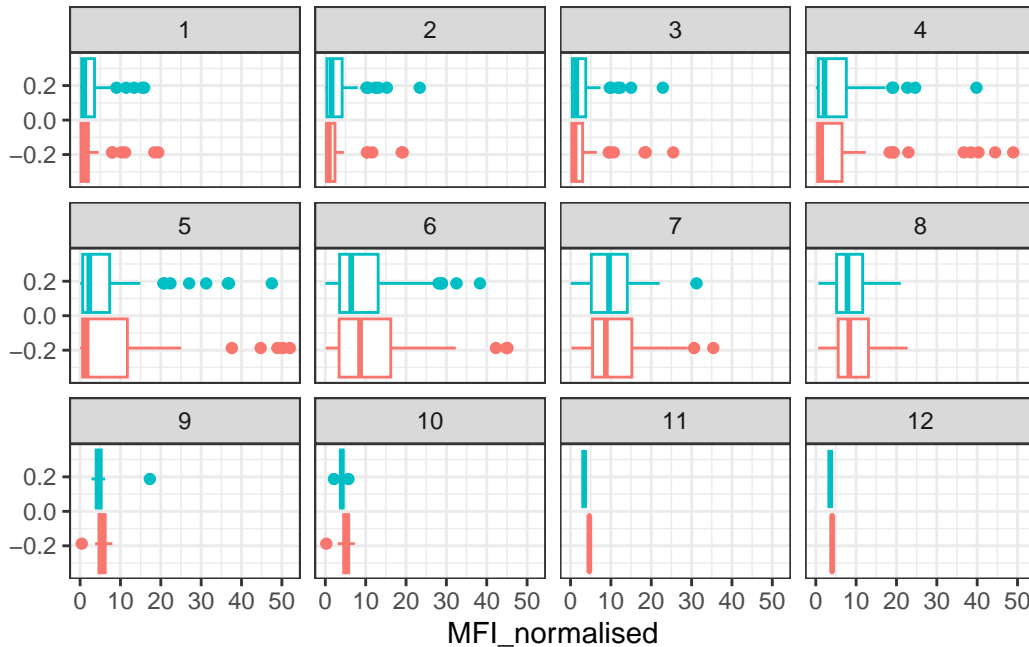
Q15. Filter to pull out only two specific antigens for analysis and create a boxplot for each. You can chose any you like. Below I picked a “control” antigen (“OVA”, that is not in our vaccines) and a clear antigen of interest (“PT”, Pertussis Toxin, one of the key virulence factors produced by the bacterium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

ANSWER:

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



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



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

ANSWER: PT antigen levels per visit (aP red, wP teal) is much higher than the OVA antigen levels per visit (aP red, wP teal). We see that at visit #5, the peak is reached then decreases for both. Both wP and aP subjects show that they peak at visit #5, and decrease.

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

ANSWER: It seems that for OVA antigen levels per visit, aP red is notably always higher than wP teal. Then, for PT antigen levels per visit, wP teal is slightly higher than aP red for each visit, until the decline.

Digging in further to look at the time course of IgG isotype PT antigen levels across aP and wP individuals:

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

ANSWER: This trend looks similar to the 2020 and 2021 dataset because they peak around the same time and decrease. However, the planned day relative to boost goes on much longer in the 2020 dataset for the wP in some cases. We can see that in 2021, the points stopped before 125; however in 2020, some of the days past 400. In addition, it seems that 2020 peaks at 7 instead of 2021 peaking at 14.

```

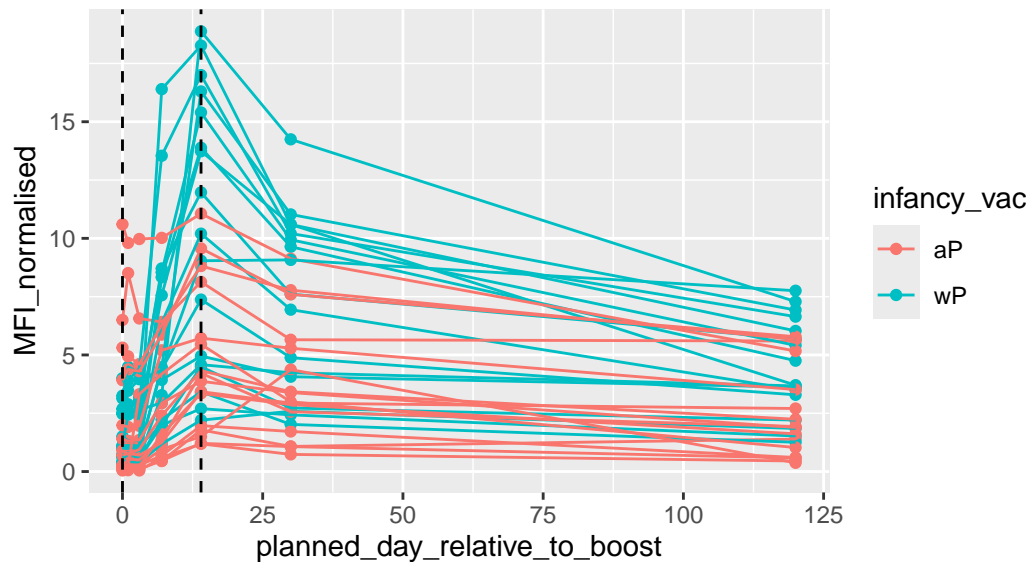
abdata.21 <- abdata %>% filter(dataset == "2021_dataset")

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)



```

## Filter to include 2021 data only
abdata.21 <- abdata %>%
  filter(dataset == "2021_dataset")

## Filter to look at IgG PT data only
pt.igg <- abdata.21 %>%

```

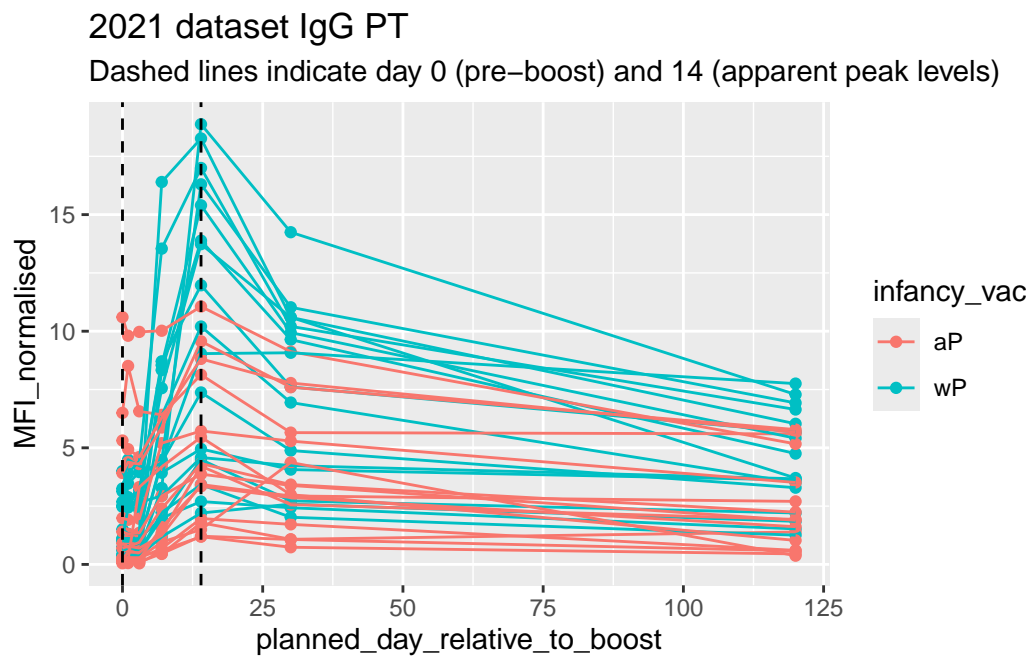


```

filter(isotype == "IgG", antigen == "PT")

## Plot and color by infancy_vac (wP vs aP)
ggplot(pt.igg) +
  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)")

```



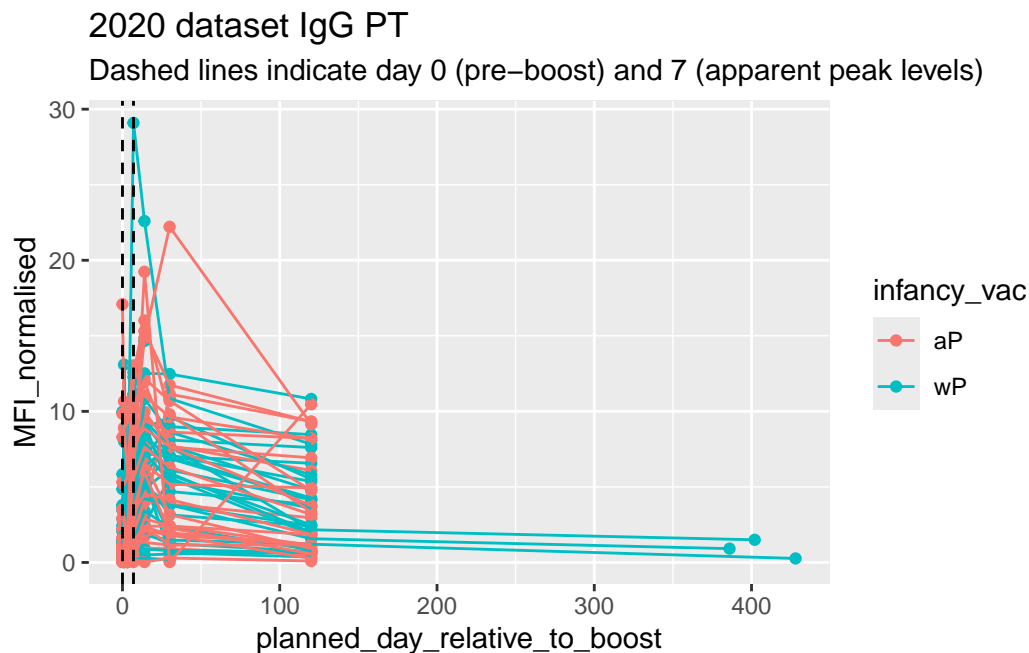
```

## Filter to include 2020 data only
abdata.20 <- abdata %>%
  filter(dataset == "2020_dataset")

## Filter to look at IgG PT data only
pt.igg <- abdata.20 %>%
  filter(isotype == "IgG", antigen == "PT")

```

```
## Plot and color by infancy_vac (wP vs aP)
ggplot(pt.igg) +
  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=7, linetype="dashed") +
  labs(title="2020 dataset IgG PT",
       subtitle = "Dashed lines indicate day 0 (pre-boost) and 7 (apparent peak levels)")
```



5. Obtaining CMI-PB RNASeq data

```
url <- "https://www.cmi-pb.org/api/v2/rnaseq?versioned_ensembl_gene_id=eq.ENSOG00000211896.7"
rna <- read_json(url, simplifyVector = TRUE)
```

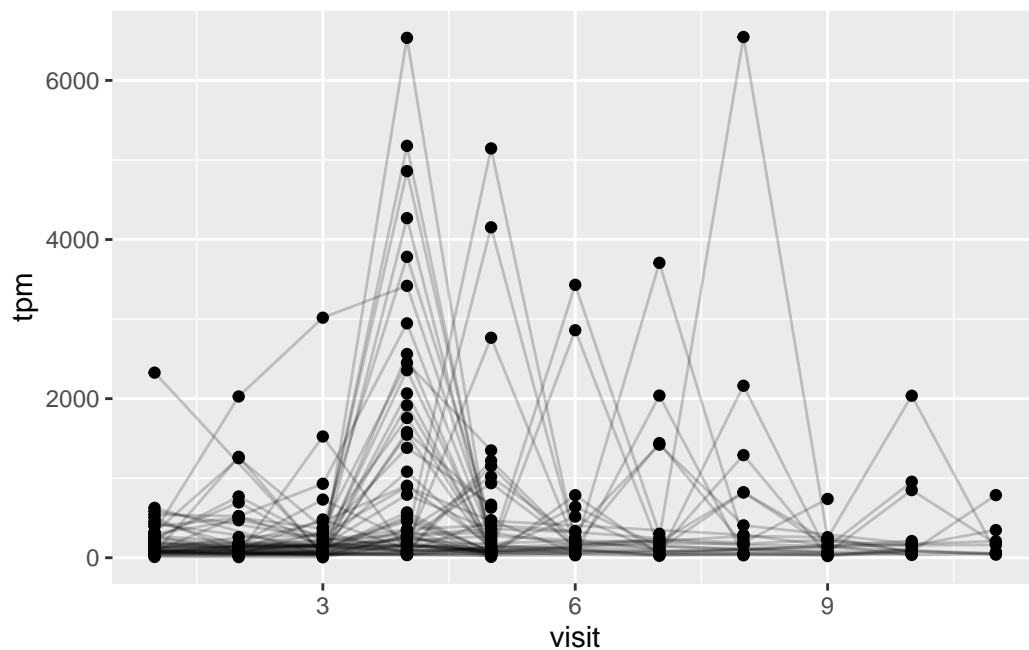
```
#meta <- inner_join(specimen, subject)
ssrna <- inner_join(rna, meta)
```

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).

ANSWER:

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



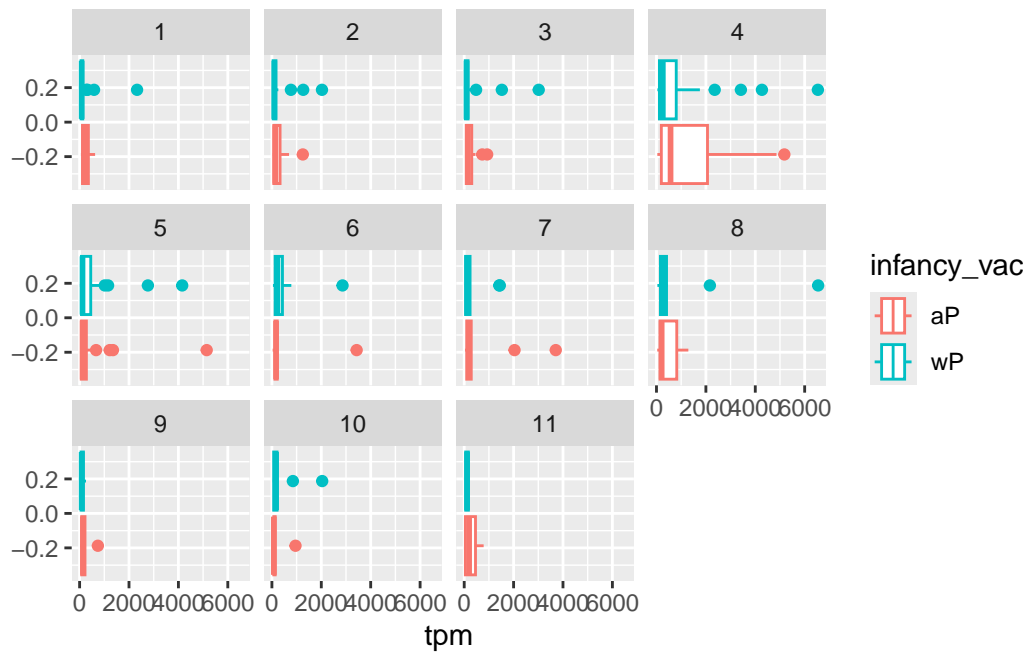
Q20.: What do you notice about the expression of this gene (i.e. when is it at it's maximum level)?

ANSWER: The expression of the gene when it is at its maximum level is around 4 visits.

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

ANSWER: This pattern in time matched the trend of the antibody titer data because they both peaked around 4-5 visits

```
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()
```

