

# ADS2 Group Exercise ICA

Group 3

14/4/2022

## Contents

<b>Load the libraries</b>	<b>1</b>
<b>Load the Dataset</b>	<b>2</b>
<b>Clean the Data</b>	<b>2</b>
Missing values . . . . .	2
Duplicates in the data . . . . .	2
Typos and Naming schemes . . . . .	3
Factoring <i>age</i> column . . . . .	4
Outliers & Strange Patterns . . . . .	4
<b>Part 1: Exploring the data</b>	<b>6</b>
Question 1 Plotting the number of deaths . . . . .	6
Question 2 Total number of malaria cases . . . . .	7
Question 3 Percentage of deaths in certain region . . . . .	9
<b>Part 2: Ask our own question</b>	<b>10</b>
Question Motivation . . . . .	10
Problem solving process . . . . .	11
Interpretion of results . . . . .	13

## Load the libraries

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

## Load the Dataset

```
malaria <- read.csv("malaria.csv")
head(malaria, 3)
```

```
##   measure      location sex      age  cause metric year
## 1 Deaths East Asia & Pacific - WB Both    Under 5 Malaria Number 2000
## 2 Deaths East Asia & Pacific - WB Both  5-14 years Malaria Number 2000
## 3 Deaths East Asia & Pacific - WB Both 15-49 years Malaria Number 2000
##      val      upper      lower
## 1 1873.7482 4692.839 741.3761
## 2  806.3444 1972.721 334.7347
## 3 4450.4837 10826.707 1905.1511
```

## Clean the Data

Firstly, we choose to explore the dataset and clean the data, since this step will benefit following analysis.

### Missing values

Deal with NA values and empty entries

```
# Change any empty entries to NA
malaria <- na_if(malaria, "")
anyNA(malaria)
```

```
## [1] FALSE
```

```
# Remove NA if have
if(anyNA(malaria)){
  malaria <- na.omit(malaria)
}
```

There are no missing values in this dataset. So we do not need to remove lines. Any missing values will be removed if future updated datasets have some.

### Duplicates in the data

Check for duplicates

```
anyDuplicated(malaria)
```

```
## [1] 0
```

```
# Remove duplicates if have
if(anyDuplicated(malaria) != 0){
  malaria <- unique(malaria)
}
```

There are no duplicated values in this dataset. So we do not need to remove lines. Any duplicated lines will be deleted if future updated datasets have some.

## Typos and Naming schemes

Test the consistency of naming scheme of some columns and whether there are some typos

See the structure of the whole dataset first

```
str(malaria)
```

```
## 'data.frame':    700 obs. of  10 variables:
##  $ measure : chr  "Deaths" "Deaths" "Deaths" "Deaths" ...
##  $ location: chr  "East Asia & Pacific - WB" "East Asia & Pacific - WB" "East Asia & Pacific - WB" "
##  $ sex      : chr  "Both" "Both" "Both" "Both" ...
##  $ age      : chr  "Under 5" "5-14 years" "15-49 years" "50 to 74 years" ...
##  $ cause    : chr  "Malaria" "Malaria" "Malaria" "Malaria" ...
##  $ metric   : chr  "Number" "Number" "Number" "Number" ...
##  $ year     : int   2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 ...
##  $ val      : num   1873.7 806.3 4450.5 2057 68.3 ...
##  $ upper    : num   4693 1973 10827 4932 161 ...
##  $ lower    : num   741.4 334.7 1905.2 872.4 29.9 ...
```

```
# location column
```

```
table(malaria$location)
```

```
##
##      East Asia & Pacific - WB      Europe & Central Asia - WB
##                100                100
## Latin America & Caribbean - WB Middle East & North Africa - WB
##                100                100
##                North America                South Asia - WB
##                100                100
##      Sub-Saharan Africa - WB
##                100
```

```
malaria$location <- gsub(' - WB', '', malaria$location) # delete ' - WB' part
table(malaria$location)
```

```
##
##      East Asia & Pacific      Europe & Central Asia
##                100                100
## Latin America & Caribbean Middle East & North Africa
##                100                100
##                North America                South Asia
##                100                100
##      Sub-Saharan Africa
##                100
```

```
# age column
```

```
table(malaria$age)
```

```
##
##      15-49 years      5-14 years 50 to 74 years      75 plus      Under 5
##                140                140                140                140                140
```

```
malaria[malaria$age == '50 to 74 years', 'age'] <- '50-74 years' # replace 'to' with '-'
table(malaria$age)
```

```
##
## 15-49 years  5-14 years 50-74 years    75 plus    Under 5
##          140      140      140        140      140
```

## Factoring *age* column

After checking the whole structure of this dataset, we want to turn *age* column into factor, in order to benefit the following plotting step.

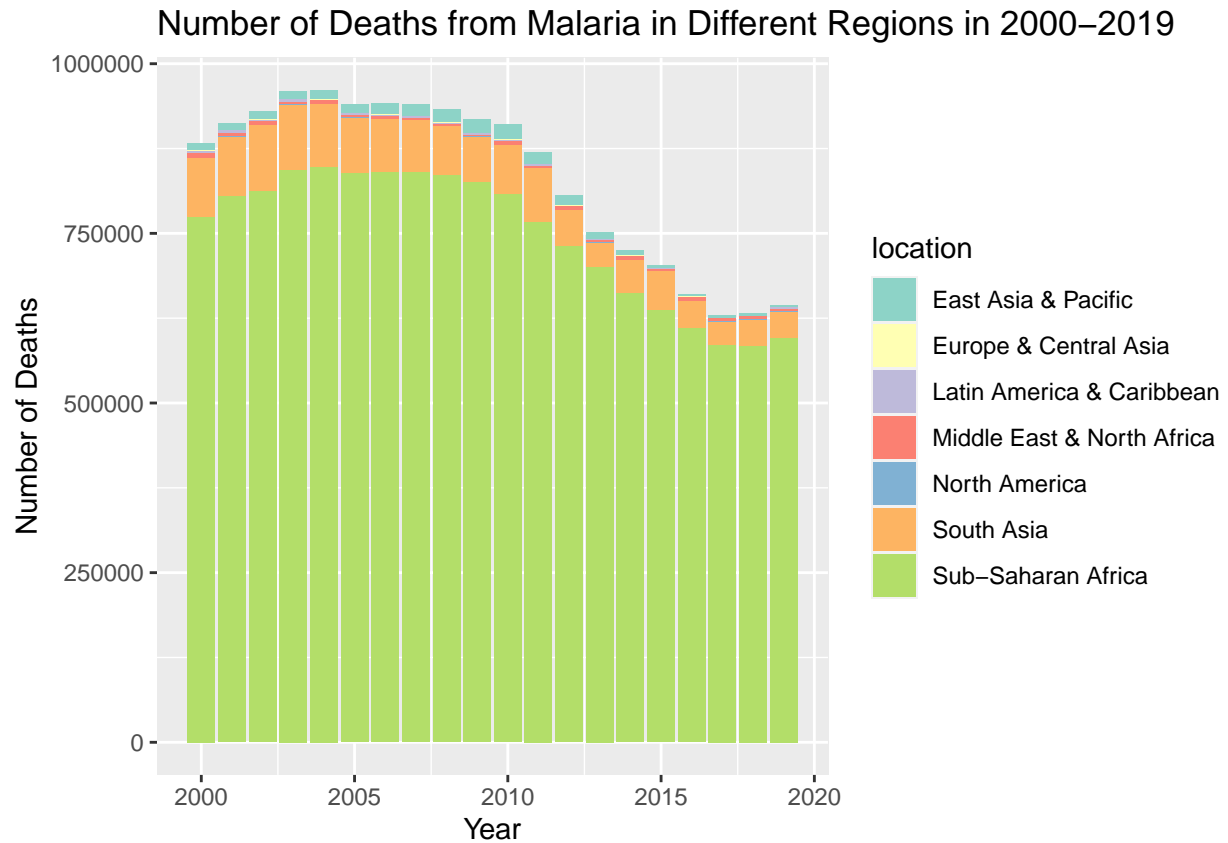
```
malaria <- malaria %>% mutate(age = factor(age,
      levels = c('Under 5', '5-14 years', '15-49 years', '50-74 years', '75 plus')))
```

## Outliers & Strange Patterns

We wonder whether the data has outstanding outliers or strange distribution patterns. So we decide to plot the number of cases against different factors to explore and examine.

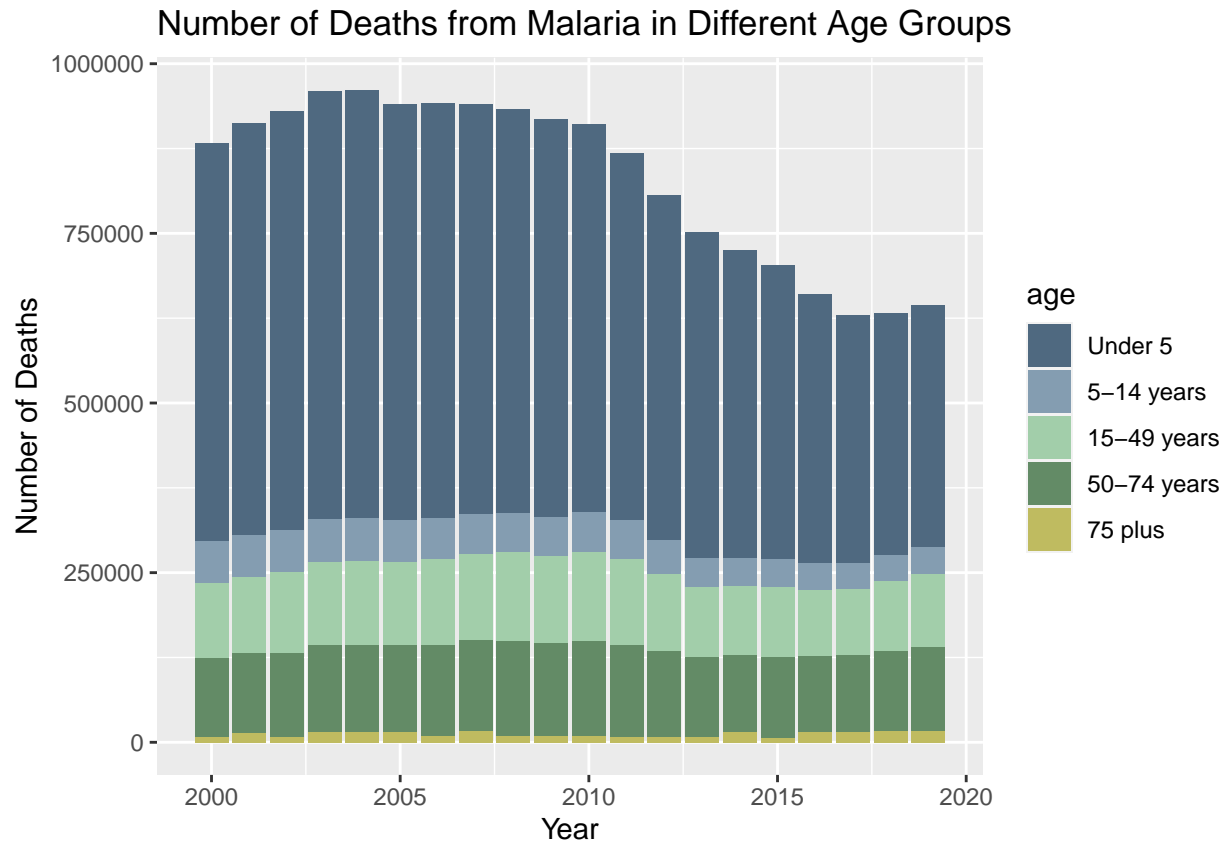
```
# deaths vs location
min_year <- min(malaria$year)
max_year <- max(malaria$year)

ggplot(malaria) +
  geom_bar(stat='identity', aes(x=year, y=val, fill=location)) +
  scale_fill_paletteer_d("RColorBrewer::Set3") +
  xlab("Year") +
  ylab("Number of Deaths") +
  ggtitle(paste("Number of Deaths from Malaria ", "in Different Regions in ",
    min_year, "-", max_year, sep = ' '))
```



We can see that the great majority of deaths of malaria happened in Sub-Saharan Africa, followed by South Asia. However, with no additional information about the malaria cases, we decided to leave the data unchanged.

```
# deaths vs age
ggplot(malaria) +
  geom_bar(stat='identity', aes(x=year, y=val, fill=age)) +
  scale_fill_paletteer_d("ggthemes::Miller_Stone") +
  xlab("Year") +
  ylab("Number of Deaths") +
  ggtitle("Number of Deaths from Malaria in Different Age Groups")
```



We can see that the majority of deaths from malaria happened in age group “Under 5”. However, with no additional information about the malaria cases, we decide to leave the data unchanged.

## Part 1: Exploring the data

In this part, we will answer the questions given in the guidance.

### Question 1 Plotting the number of deaths

**Question:** Plot the number of deaths from malaria between 2000 and 2019 for each of the age groups for the East Asia and Pacific region. What age group seems to have the highest number of cases, and why do you think that is?

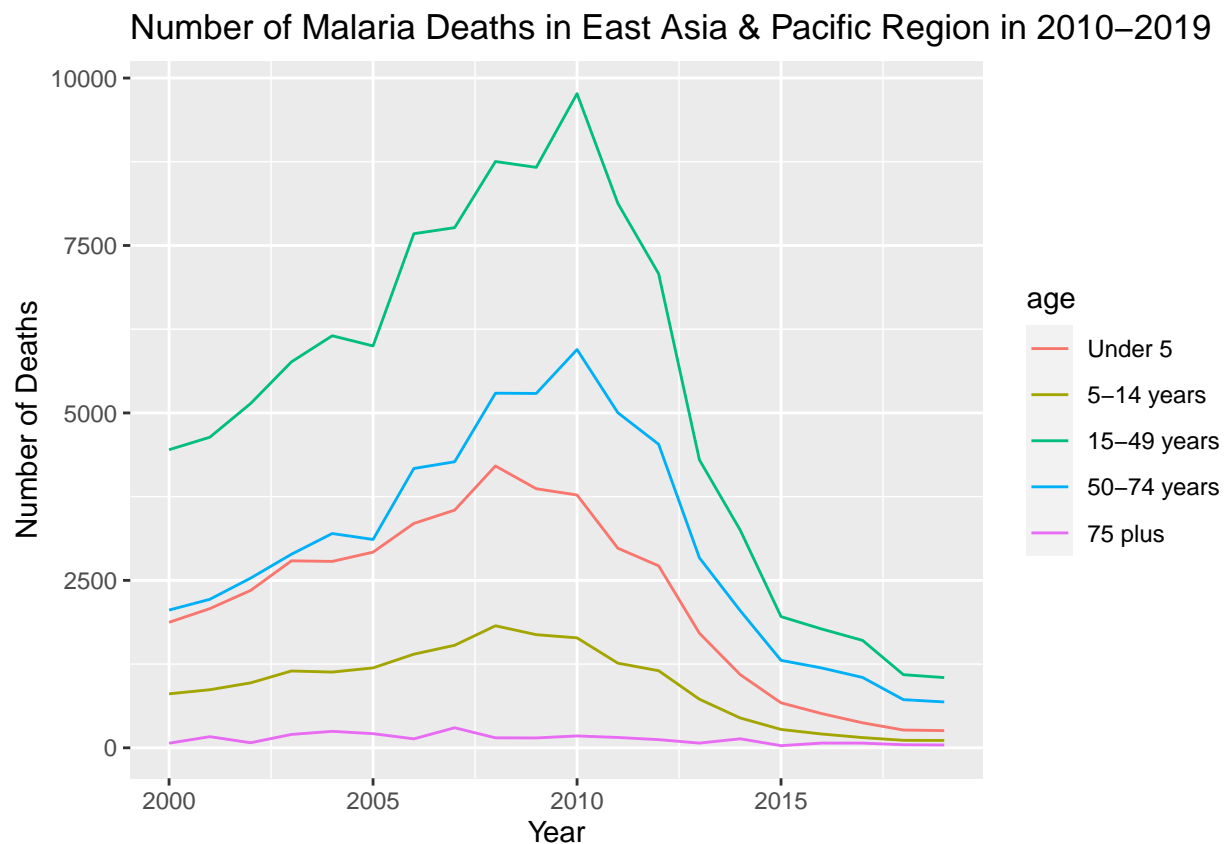
To solve this problem, firstly we need to extract data of the East Asia and Pacific region from the original dataset.

```
subregion <- malaria %>%
  filter(location == 'East Asia & Pacific') %>%
  mutate(age=factor(age, levels=c('Under 5','5-14 years',
                                  '15-49 years','50-74 years','75 plus')))
head(subregion, 2)
```

```
##   measure      location sex      age  cause metric year      val
## 1 Deaths East Asia & Pacific Both    Under 5 Malaria Number 2000 1873.7482
## 2 Deaths East Asia & Pacific Both 5-14 years Malaria Number 2000 806.3444
##      upper      lower
## 1 4692.839 741.3761
## 2 1972.721 334.7347
```

Then we can plot the data depending on the time and age groups.

```
ggplot(subregion, aes(x = year, y = val, group = age)) +
  geom_line(aes(color = age)) +
  xlab("Year") +
  ylab("Number of Deaths") +
  ggtitle("Number of Malaria Deaths in East Asia & Pacific Region in 2010–2019")
```



As we can see from the results, it seems that the age group “15-49 years” always have the highest number of cases which died from malaria in every year between 2000 and 2019.

It may be because that this group has the largest population compared with other age groups. Also, it is possible that this age group is optimal labour force, so they are more likely to get in touch with other people, so this age group may have higher probability to get infected by malaria.

## Question 2 Total number of malaria cases

**Question:** In which year was the total number of malaria cases (across all regions and age groups) the highest? In which year was it the lowest?

To solve this problem, we can transform the data set, and summarize number of deaths grouped by year.

```
malaria_year <- malaria %>%
  group_by(year) %>%
  summarise(val=sum(val))
head(malaria_year, 3)
```

```
## # A tibble: 3 x 2
##   year    val
##   <int>  <dbl>
## 1  2000 882060.
## 2  2001 912710.
## 3  2002 929197.
```

```
malaria_max <- malaria_year[malaria_year$val == max(malaria_year$val), 'year']
malaria_min <- malaria_year[malaria_year$val == min(malaria_year$val), 'year']
```

```
# Conclusion 1
print(paste("Total number of malaria cases was highest in", malaria_max))
```

```
## [1] "Total number of malaria cases was highest in 2004"
```

```
# Conclusion 2
print(paste("Total number of malaria cases was lowest in", malaria_min))
```

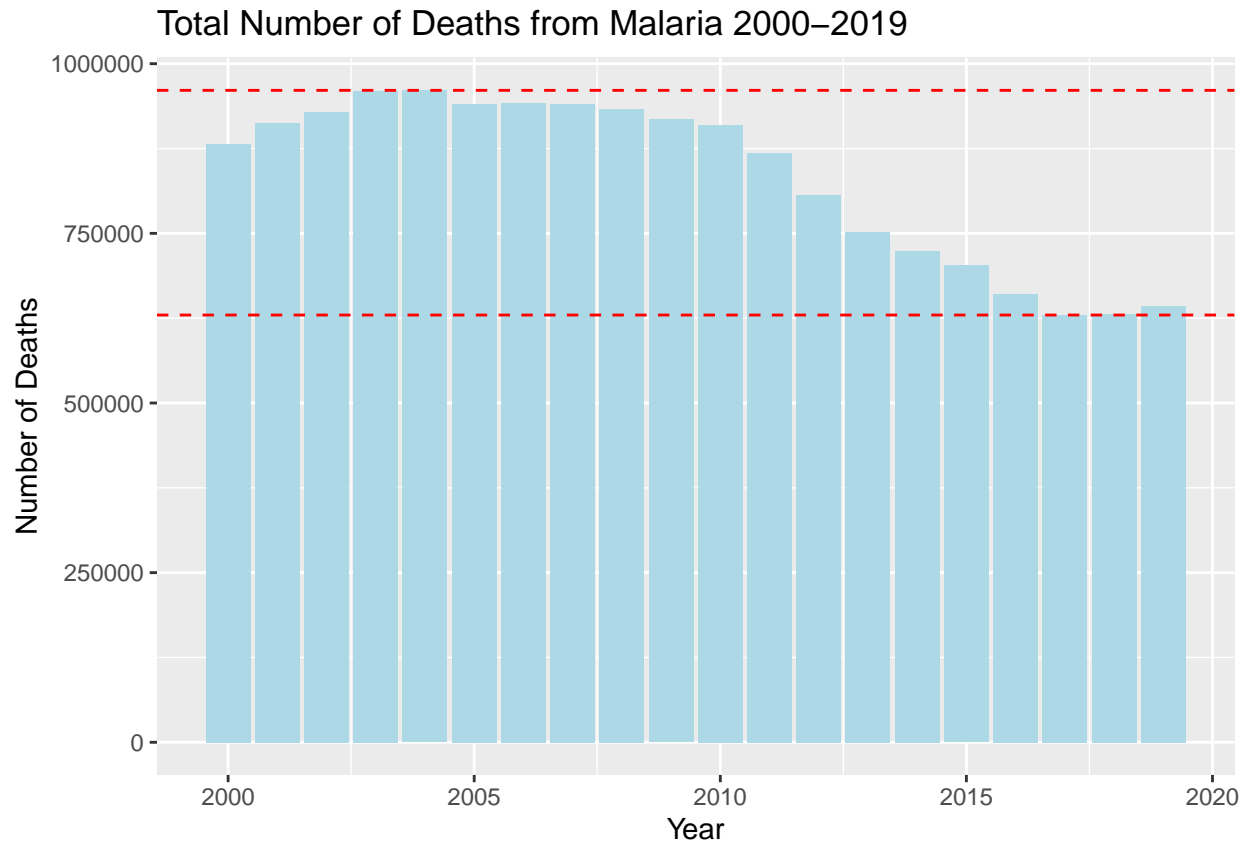
```
## [1] "Total number of malaria cases was lowest in 2017"
```

This can also be verified by visualizing the number of total death each year.

```
ymax <- malaria_year[malaria_year$year == malaria_max$year, "val"]
ymin <- malaria_year[malaria_year$year == malaria_min$year, "val"]

ggplot(malaria_year) +
  geom_bar(stat='identity', aes(x=year, y=val), fill="lightblue") +
  xlab("Year") +
  ylab("Number of Deaths") +
  ggtitle(paste("Total Number of Deaths from Malaria ",
                min_year, "-", max_year, sep = ' ')) +
  geom_hline(col = "red", linetype = "dashed", yintercept = ymax$val) +
  geom_hline(col = "red", linetype = "dashed", yintercept = ymin$val)
```





### Question 3 Percentage of deaths in certain region

**Question:** What percentage of total Malaria deaths in 2010 happened in the Latin America and Caribbean region?

To solve this problem, we can transform the data set and summarize number of deaths in 2010 grouped by location.

```
malaria_location <- malaria %>%
  filter(year == "2010") %>%
  group_by(location) %>%
  summarise(val=sum(val))
head(malaria_location, 2)
```

```
## # A tibble: 2 x 2
##   location          val
##   <chr>            <dbl>
## 1 East Asia & Pacific 21306.
## 2 Europe & Central Asia 0.936
```

Then we can choose total Malaria deaths happened in the Latin America and Caribbean region.

```

malaria_latin_cari <- malaria_location %>%
  filter(location=='Latin America & Caribbean')
malaria_latin_cari

```

```

## # A tibble: 1 x 2
##   location          val
##   <chr>            <dbl>
## 1 Latin America & Caribbean 2675.

```

```

total_val <- sum(malaria_location$val)
total_val

```

```
## [1] 909816.9
```

```

# Conclusion
print(paste("About ", round(100*malaria_latin_cari$val/total_val, digits = 2),
  "% of total Malaria deaths in 2010 happened ",
  "in the Latin America and Carribean region.",
  sep = ""))

```

```
## [1] "About 0.29% of total Malaria deaths in 2010 happened in the Latin America and Carribean region."
```

## Part 2: Ask our own question

In this part, we will ask one question that we are interested in, and choose a suitable method to use the data provided to answer it.

### Question Motivation

From the previous questions and visualizations, we can see that the number of malaria deaths has been constantly changing each year. Thus, our group wants to investigate how the number of malaria deaths changes over years in a specific region, for example, **in Middle East and North Africa**.

This can reflect the effectiveness of malaria control in the area at a given time. That is, if we can see a significant decrease in malaria deaths between two years, this could mean that the approach to control malaria in this area has been effective over this period of time.

Before using any statistical tests, we first plot the percentage of deaths each year in Middle East & North Africa region.

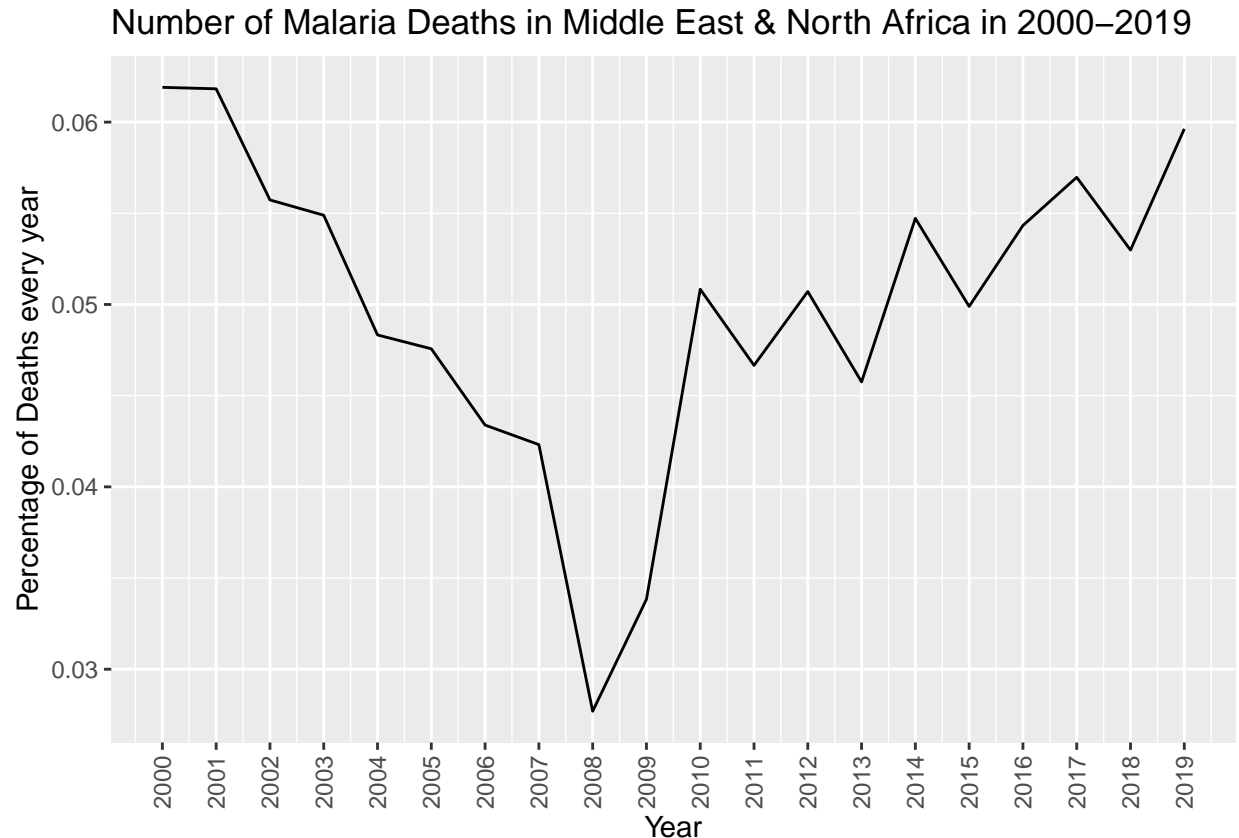
```

# Malaria deaths (integer) in Middle East & North Africa
malaria_ME_NA <- malaria %>%
  filter(location == 'Middle East & North Africa') %>%
  group_by(year) %>%
  summarise(val=sum(val))

malaria_ME_NA$val <- as.integer(malaria_ME_NA$val) # change to integer format
total_death <- sum(malaria_ME_NA$val) # total death in this region over the years

```

```
# plot the percentage of deaths each year
ggplot(malaria_ME_NA) +
  geom_line(stat = 'identity', aes(x=year, y=val/total_death)) +
  xlab('Year') +
  ylab('Percentage of Deaths every year') +
  ggtitle(paste("Number of Malaria Deaths in Middle East & North Africa in ",
                min_year, "-", max_year, sep = ' ')) +
  scale_x_continuous(breaks = seq(min_year, max_year, 1)) +
  theme(axis.text.x = element_text(angle=90, hjust=1, vjust=0.5))
```



From the plot we can roughly observe a decreasing trend of the percentages of deaths from 2000-2008.

But wait, is there really a difference between, for example, percentages of death in 2006 and 2007, or is this variation likely to be caused by chance? What about between 2012 and 2013?

**Question:** Is there a significant difference in the percentage of people dying from malaria each year in the Middle East and North Africa?

This percentage refers to the number of deaths of that region in that year / the total number of deaths of that region in all years.

## Problem solving process

Since we want to test whether there is a real difference between two samples, we have to formulate our hypothesis first.

**Null Hypothesis:** There is no significant difference of percentage of malaria deaths between any two years in Middle East & North Africa region.

**Alternate Hypothesis:** The differences of percentage of malaria deaths in Middle East & North Africa region are significant at least in two years.

**Step1:** To address this question, we firstly explore whether any statistical methods like t-test or ANOVA can be used. After exploring the data, we find that there is only one value for each year, representing the number of malaria deaths that year. That is, we do not know the population distribution. Finally, we decide to adopt a bootstrap test.

The type of bootstrapping method we choose is *Case Resampling* approach.

The parameter we want to compare in bootstrapping is the percentage of malaria deaths.

To make the code of bootstrapping more concise, we simplify the calculation process of the percentage. When generating mock samples of each year, we use 1 to represent deaths in that year and 0 for deaths in other years. After sampling with replacement, the percentage can be regarded as the mean of annual sample, because when we calculate the average, we always have to divide by the sum.

```
# Bootstrapping for each year
quantiles.total <- c()
for (y in malaria_ME_NA$year) {
  current_death <- malaria_ME_NA %>%
    filter(year == y) %>%
    select(val) %>%
    .$val

  boot.year <- c(rep(1, current_death), rep(0, total_death-current_death))
  year.boot.300 <- replicate(300, mean(sample(boot.year, size = total_death, replace = T)))
  quantiles.total <- c(quantiles.total, quantile(year.boot.300, c(0.025, 0.975)))
}
head(quantiles.total, 5)
```

```
##          2.5%          97.5%          2.5%          97.5%          2.5%
## 0.06034316 0.06373369 0.06012531 0.06348393 0.05412801
```

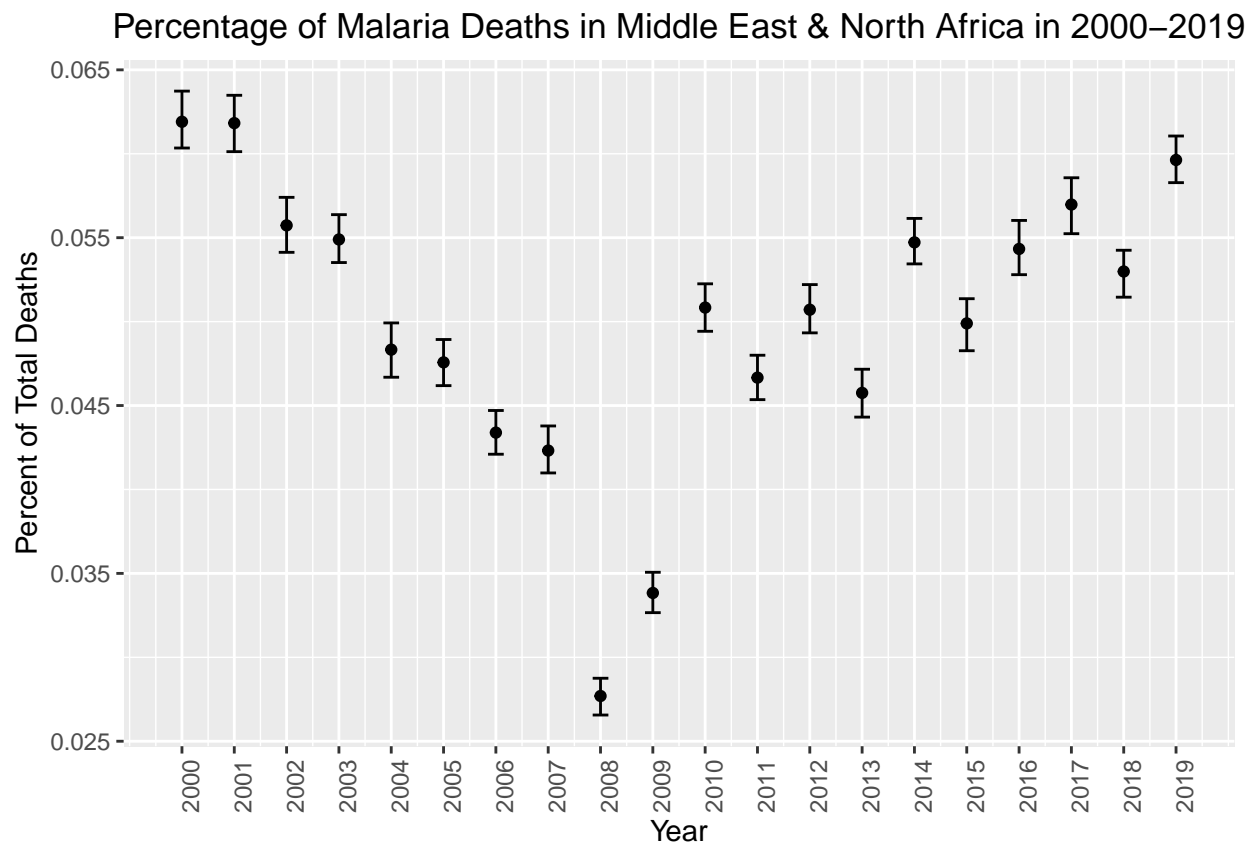
```
# Write down the bootstrapping results
matrix.quantile <- matrix(data = quantiles.total, nrow = length(quantiles.total)/2,
                          ncol = 2, byrow = T)

df.quantile <- as.data.frame(matrix.quantile)
df.quantile$year <- malaria_ME_NA$year
df.quantile$perc <- malaria_ME_NA$val/total_death
head(df.quantile)
```

```
##          V1          V2 year      perc
## 1 0.06034316 0.06373369 2000 0.06190723
## 2 0.06012531 0.06348393 2001 0.06182675
## 3 0.05412801 0.05740358 2002 0.05573375
## 4 0.05351986 0.05637265 2003 0.05489452
## 5 0.04668564 0.04992326 2004 0.04833017
## 6 0.04618526 0.04893401 2005 0.04757142
```

**Step2:** Visualize the 95% confidence interval data to check the differences between groups

```
# Visualization
ggplot(df.quantile) +
  geom_point(aes(x=year, group=year, y=perc)) +
  geom_errorbar(aes(x=year, ymin=V1, ymax=V2), width=0.3) +
  xlab('Year') +
  ylab('Percent of Total Deaths') +
  ggtitle(paste("Percentage of Malaria Deaths in Middle East & North Africa in ",
    min_year, "-", max_year, sep = ' ')) +
  scale_x_continuous(breaks = seq(min_year, max_year, 1)) +
  theme(plot.title = element_text(hjust=0.5),
    axis.text.x = element_text(angle=90, hjust=1))
```



## Interpretion of results

As we can see from the plot above, we have sufficient evidence to reject  $H_0$ . That is, the differences in percentage deaths from malaria in the Middle East & North Africa region are significant **between some two years**.

We would expect the 95% confidence intervals drawn above to be non-overlapping for groups significantly different from each other, such as year 2012 and 2013. Therefore, there is a significant difference for the percentage of Malaria deaths between 2012 and 2013. However, there is no significant difference between 2006 and 2007 since they have overlaps during the 95% confidence intervals.

More generally, based on the same principle, we can analyze if there is a significant change in malaria deaths between any years in a given region. This change is a crucial indicator for evaluating whether the malaria situation in that region has improved, worsened or remained unchanged between the two years.