

MERS-CoV Korea 2015 - Marta Smith

2025-12-01

Introduction

In this R Markdown, I analyze the `mers_korea_2015` dataset from the “outbreaks” R package and answer questions about the MERS-CoV outbreak in South Korea in 2015. Source: <https://github.com/luisfmontemayor/learn-compbio/tree/main>

Loading and inspecting the data

```
# loading the data  
data("mers_korea_2015")
```

Quick questions

1. What object class is `mers_korea_2015` in R? Describe what this class is like, and how is it different from a simple atomic type in R (like a `string` like “hello world!” or a `logical` value like `true`) or a vector type defined with `c()`

`mers_korea_2015` is a list. A list is a collection of data which is ordered and changeable (source).

2. What sub-items exist within this item?

It contains 2 sub-items that are dataframes: `linelist` and `contacts`. `linelist`: A dataframe of MERS-CoV cases and their attributes (source). `contacts`: A dataframe describing the relationship between MERS-CoV cases (source).

```
# inspecting the raw data  
class(mers_korea_2015)
```

```
## [1] "list"
```

```
class(mers_korea_2015$linelist)
```

```
## [1] "data.frame"
```

```
class(mers_korea_2015$contacts)
```

```
## [1] "data.frame"
```

```
length(mers_korea_2015) # number of elements the list contains
```

```
## [1] 2
```

```
names(mers_korea_2015) # names of the elements of the list
```

```
## [1] "linelist" "contacts"
```

```
str(mers_korea_2015) # structure function
```

```
## List of 2
## $ linelist:'data.frame': 162 obs. of 15 variables:
## ..$ id : chr [1:162] "SK_1" "SK_2" "SK_3" "SK_4" ...
## ..$ age : int [1:162] 68 63 76 46 50 71 28 46 56 44 ...
## ..$ age_class : chr [1:162] "60-69" "60-69" "70-79" "40-49" ...
## ..$ sex : Factor w/ 2 levels "F","M": 2 1 2 1 2 2 1 1 2 2 ...
## ..$ place_infect : Factor w/ 2 levels "Middle East",...: 1 2 2 2 2 2 2 2 2 2 ...
## ..$ reporting_ctype: Factor w/ 2 levels "China","South Korea": 2 2 2 2 2 2 2 2 2 1 ...
## ..$ loc_hosp : Factor w/ 13 levels "365 Yeollin Clinic, Seoul",...: 10 10 10 10 1 10 10 13 10 ...
## ..$ dt_onset : Date[1:162], format: "2015-05-11" "2015-05-18" ...
## ..$ dt_report : Date[1:162], format: "2015-05-19" "2015-05-20" ...
## ..$ week_report : Factor w/ 5 levels "2015_21","2015_22",...: 1 1 1 2 2 2 2 2 2 2 ...
## ..$ dt_start_exp : Date[1:162], format: "2015-04-18" "2015-05-15" ...
## ..$ dt_end_exp : Date[1:162], format: "2015-05-04" "2015-05-20" ...
## ..$ dt_diag : Date[1:162], format: "2015-05-20" "2015-05-20" ...
## ..$ outcome : Factor w/ 2 levels "Alive","Dead": 1 1 2 1 1 2 1 1 1 1 ...
## ..$ dt_death : Date[1:162], format: NA NA ...
## $ contacts:'data.frame': 98 obs. of 4 variables:
## ..$ from : chr [1:98] "SK_14" "SK_14" "SK_14" "SK_14" ...
## ..$ to : chr [1:98] "SK_113" "SK_116" "SK_41" "SK_112" ...
## ..$ exposure : Factor w/ 5 levels "Contact with HCW",...: 2 2 2 2 2 2 2 2 2 2 ...
## ..$ diff_dt_onset: int [1:98] 10 13 14 14 15 15 15 16 16 16 ...
```

3. How can one access the sub-items within the `mers_korea_2015`?

```
# separating the data sets
# different ways of accessing the sub-items: either with a $ or [1]
mers_korea_2015$linelist
mers_korea_2015[1]

# creating new variables for each sub-item
linelist <- mers_korea_2015$linelist
contacts <- mers_korea_2015$contacts
```

4. The “sub-items” within `mers_korea_2015` are different data sets within our project. Play with them! Open them, plot graphs, get your hands dirty. The questions won’t always tell you what data set to use! So know what you’re starting with, to make sure that you can

```

# cleaning the data
linelist <- linelist %>%
  mutate(
    loc_hosp = str_replace(
      string = loc_hosp,
      pattern = ",.*",
      replacement = ""
    ) ) %>%
  mutate(
    loc_hosp = str_trim(loc_hosp)
  )
unique(linelist$loc_hosp)

```

```

## [1] "Pyeongtaek St. Mary"
## [2] "365 Yeollin Clinic"
## [3] "Seoul Clinic"
## [4] "Konyang University Hospital"
## [5] "Dae Cheong Hospital"
## [6] "Samsung Medical Center"
## [7] "Hallym University Medical Center"
## [8] "Hallym University Dongtan Sacred Heart Hospital"
## [9] "Pyeongtaek goodmorning hospital"
## [10] "Pyeongtaek fraternity hospital"
## [11] "Konkuk University hospital"
## [12] "Gangdong Kyuang Hee University Hospital"

```

Analyses

Getting my hands dirty with the data

```

# admissions grouped by sex and date of report
admissions_summary <- linelist %>%
  group_by(dt_report, sex) %>%
  summarise(
    count = n(),
    mean_age = mean(age, na.rm= TRUE),
    .groups = 'drop'
  )

# Pivot table by ages per date of report
pivot_age <- admissions_summary %>%
  pivot_wider(
    id_cols = dt_report, # columns
    names_from = sex, # names of the column
    values_from = mean_age, # values
    names_prefix = "mean_age_" # optional
  )

# count of admissions by the number of week
admissions_by_week <- ggplot(linelist) +

```

```

geom_bar(aes(x=week_report), fill="gray", color="black") + # color is used for the outline
theme_minimal() +
labs(
  title="MERS admissions by week",
  x="Week",
  y="Admissions"
)

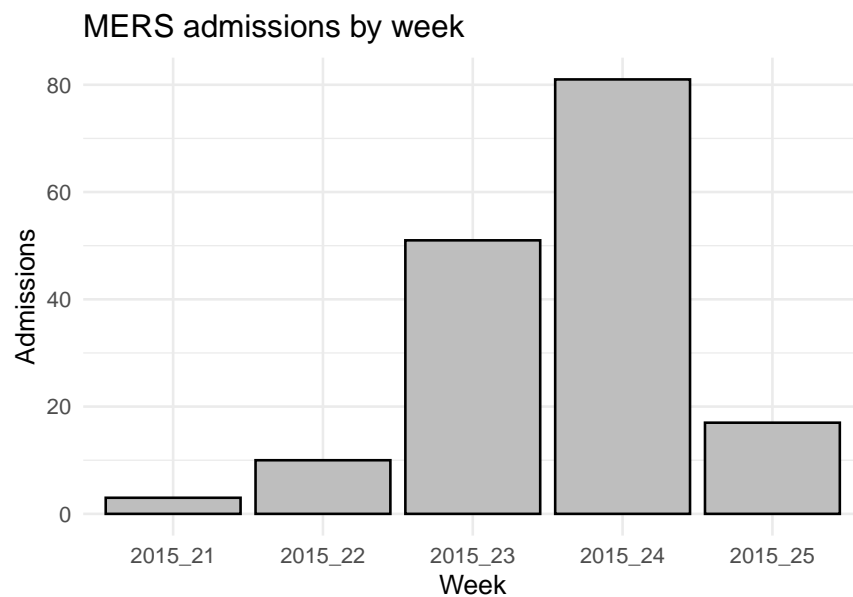
# admissions by date of report
admissions_by_dt_report <- ggplot(linelist) +
  geom_bar(aes(x=dt_report), fill="gray", color="black") + # color is used for the outline
  theme_minimal() +
  labs(
    title="MERS admissions by date of report",
    x="Date of Report",
    y="Admissions"
  )

# admissions by sex
admissions_by_sex <- ggplot(admissions_summary, aes(x=dt_report, y=count, color=sex)) +
  geom_line(linewidth =0.7, alpha=0.7) + geom_point(size=2, alpha=0.7)

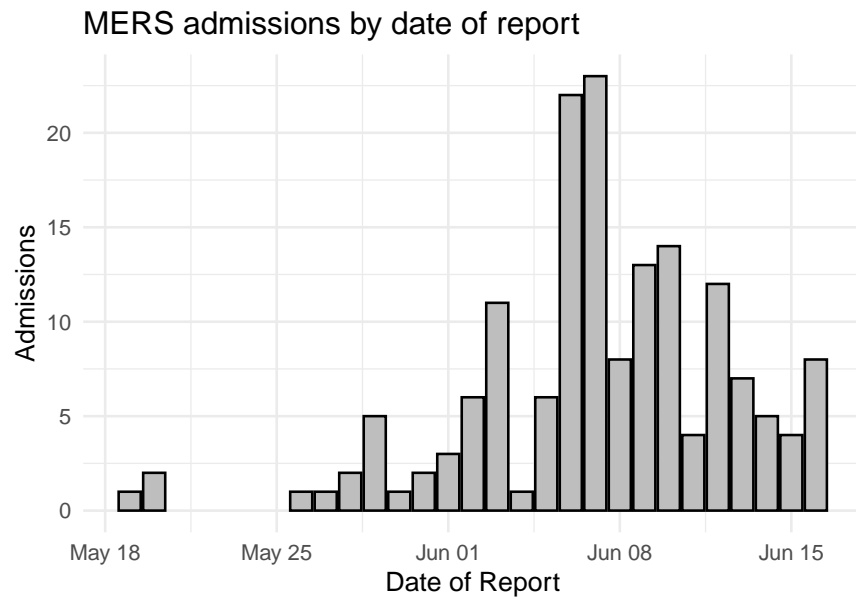
# first approach to understand the average age of the two populations
age_by_outcome <- ggplot(linelist) +
  geom_violin(aes(x=outcome, y=age, fill=outcome)) +
  theme_minimal() +
  labs(title="",
    x="Outcome", y="Age", fill="Outcome")

```

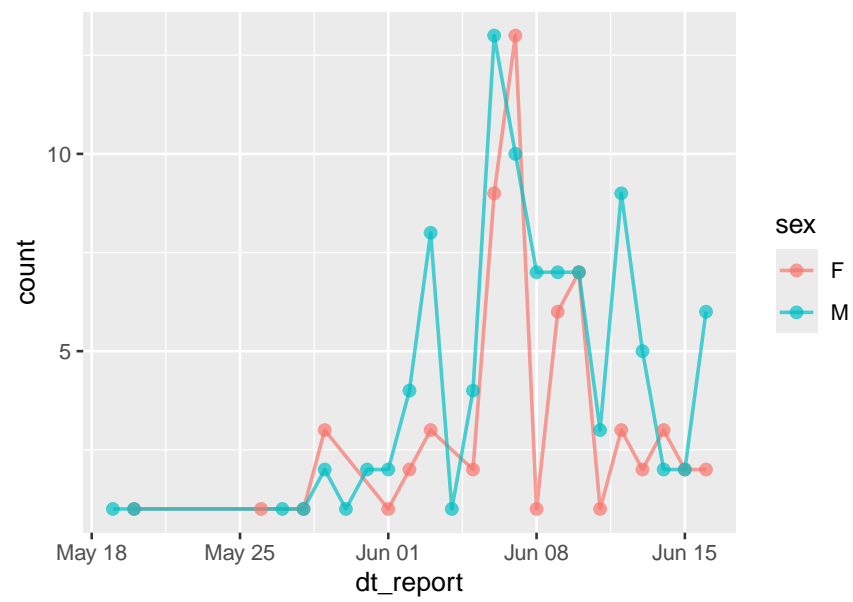
admissions_by_week



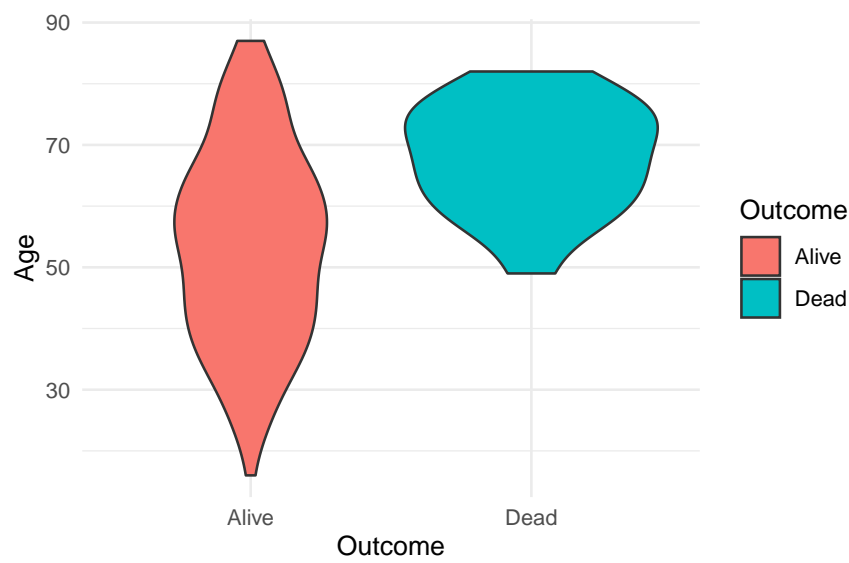
admissions_by_dt_report



admissions_by_sex



age_by_outcome



1. Demographics

1. Calculate the average age for patients with the outcome “Dead” and patients with the outcome “Alive.”

Average ages: Dead: 68 years old; Alive: 53.6 years old. The average age of this sample was 55 years old.

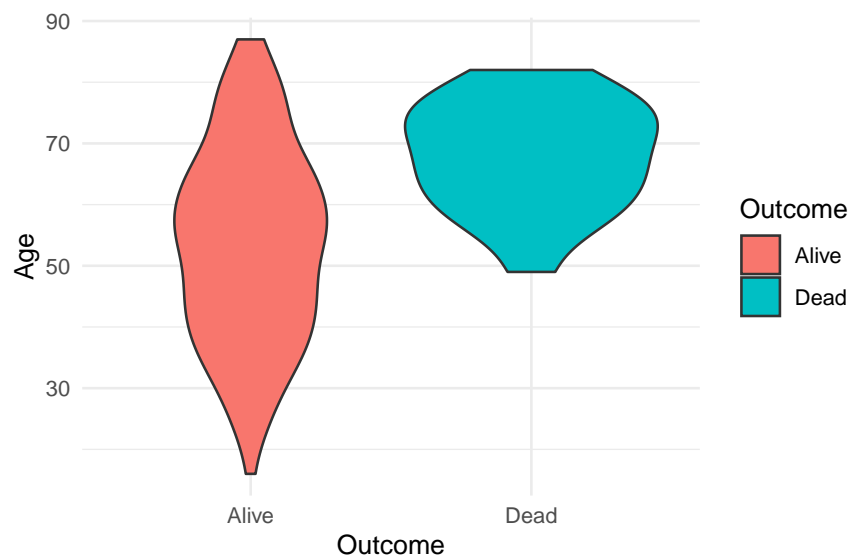
```
# age distribution (not grouped)
age_stats <- linelist %>%
  summarise(
    count = n(),
    mean_age = mean(age, na.rm= TRUE),
    median_age = median(age, na.rm = TRUE),
    sd_age = sd(age, na.rm = TRUE),
    .groups = 'drop'
  )
age_stats
```

```
##   count mean_age median_age sd_age
## 1    162 55.32099         56 15.814
```

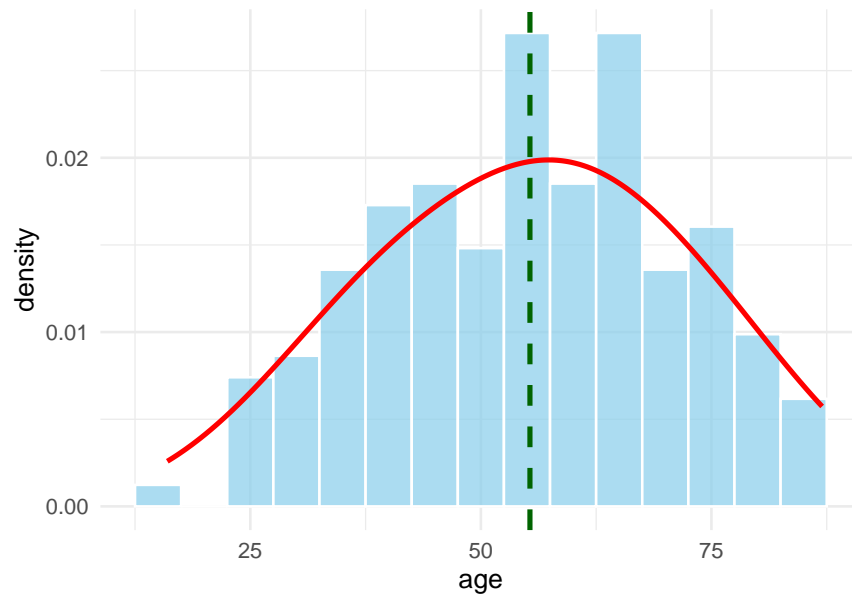
```
# summary table (grouped by outcome)
outcome_stats <- linelist %>%
  group_by(outcome) %>%
  summarise(
    count = n(),
    mean_age = mean(age, na.rm= TRUE),
    median_age = median(age, na.rm = TRUE),
    sd_age = sd(age, na.rm = TRUE),
    .groups = 'drop'
  )
outcome_stats
```

```
## # A tibble: 2 x 5
##   outcome count mean_age median_age sd_age
##   <fct>   <int>   <dbl>     <int> <dbl>
## 1 Alive    143    53.6         55  15.8
## 2 Dead     19    68.1         68   8.90
```

```
# ggplot for the average age of the two populations
age_by_outcome <- ggplot(linelist) +
  geom_violin(aes(x=outcome, y=age, fill=outcome)) +
  theme_minimal() +
  labs(title="",
        x="Outcome", y="Age", fill="Outcome")
age_by_outcome
```



```
# Age histogram + normal distribution + mean age
age_histogram <- ggplot(linelist) +
  geom_histogram(
    aes(x = age, y = after_stat(density)), # y = density for overlay
    binwidth = 5, fill = "skyblue", color = "white", alpha = 0.7 ) +
  geom_density(
    aes(x = age),
    color = "red",
    linewidth = 1,
    adjust = 2
  ) +
  geom_vline( # Average age as vertical line
    xintercept = age_stats$mean_age,
    color = "darkgreen",
    linetype = "dashed",
    linewidth = 1 ) +
  theme_minimal()
age_histogram
```



2. Based on first impressions, does age seem to be a determinant of mortality in this outbreak? How might the social status of the elderly in South Korea contribute to exposure risks (e.g., care homes, hospital frequency)?

In South Korea, the elderly population is comprised of individuals aged 65 and older (https://en.wikipedia.org/wiki/Aging_of_South_Korea). In this dataset, 49 of the cases are ages 65 and older (30% of the dataset).

Age seemed to be a determinant of mortality during the MERS-CoV outbreak in South Korea. The deceased group represents 11.7% of the dataset. The average age of the Dead group is 68 years old, 14.4 years higher than the Alive group (53.6). Also, the standard deviation for the Dead group is 56% lower in the Dead group with respect to the Alive group (8.9 vs 15.8). This shows that the data points in the Dead group are closer to the mean.

Analyzing the types of exposures in this outbreak, 47% of all infections came from Hospital Room visits (27%) and Emergency Rooms (20%). “The index case was a returning traveler from the Middle East. The infection had spread within the hospital, and subsequently to other hospitals because of patient movement, resulting in nosocomial transmission at 16 clinics and hospitals.” (Source: <https://pmc.ncbi.nlm.nih.gov/articles/PMC5840604/>).

```
hospital_count <- linelist %>%
  group_by(loc_hosp) %>%
  summarise(
    count = n(),
    percentage = count / nrow(linelist) * 100,
    mean_age = mean(age, na.rm = TRUE),
    .groups = 'drop'
  ) %>%
  arrange(desc(count))
head(hospital_count)
```

```
## # A tibble: 6 x 4
##   loc_hosp                count percentage mean_age
##   <chr>                  <int>      <dbl>    <dbl>
## 1 Samsung Medical Center      85      52.5     53.9
```


## 2	Pyeongtaek St. Mary	37	22.8	51.7
## 3	Dae Cheong Hospital	12	7.41	67.1
## 4	Konyang University Hospital	11	6.79	69.6
## 5	Hallym University Medical Center	4	2.47	59.5
## 6	Pyeongtaek goodmorning hospital	4	2.47	73.2

```

exposure_count <- contacts %>%
  group_by(exposure) %>%
  summarise(
    count = n(),
    percentage = count / nrow(linelist) * 100,
    .groups = 'drop'
  ) %>%
  arrange(desc(count))
head(exposure_count)

```

```

## # A tibble: 5 x 3
##   exposure      count percentage
##   <fct>         <int>      <dbl>
## 1 "Hospital room"      44      27.2
## 2 "Emergency room"    33      20.4
## 3 "Visit hospital"   12       7.41
## 4 "Contact with HCW"  8       4.94
## 5 "Family member "    1       0.617

```

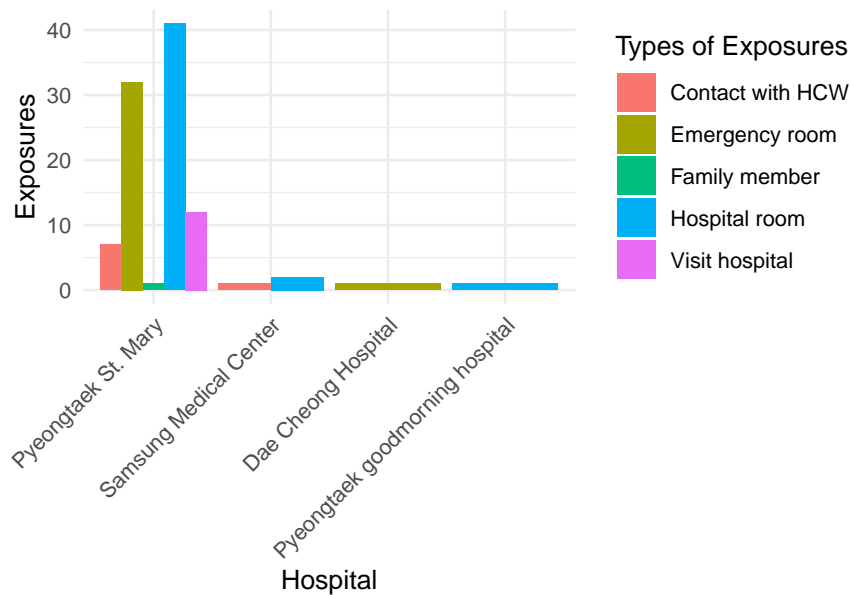
```

hospital_dist <- contacts %>%
  rename(id = from) %>%
  select(id, exposure) %>%
  left_join(linelist, by="id") %>%
  group_by(loc_hosp, exposure) %>%
  summarise(
    count = n(),
    percentage = count / nrow(linelist) * 100,
    mean_age = mean(age, na.rm=TRUE),
    .groups = 'drop'
  ) %>%
  arrange(desc(count))
(hospital_dist)

```

##	loc_hosp	exposure	count	percentage	mean_age
##	<chr>	<fct>	<int>	<dbl>	<dbl>
## 1	Pyeongtaek St. Mary	"Hospital room"	41	25.3	47.0
## 2	Pyeongtaek St. Mary	"Emergency room"	32	19.8	36.3
## 3	Pyeongtaek St. Mary	"Visit hospital"	12	7.41	68
## 4	Pyeongtaek St. Mary	"Contact with HCW"	7	4.32	53.9
## 5	Samsung Medical Center	"Hospital room"	2	1.23	70
## 6	Dae Cheong Hospital	"Emergency room"	1	0.617	78
## 7	Pyeongtaek St. Mary	"Family member "	1	0.617	68
## 8	Pyeongtaek goodmorning hospital	"Hospital room"	1	0.617	67
## 9	Samsung Medical Center	"Contact with HCW"	1	0.617	75

```
ggplot(hospital_dist, aes(x = fct_infreq(loc_hosp), y = count, fill=exposure)) +
  geom_bar(stat="identity", position="dodge") +
  labs (
    x="Hospital",
    y="Exposures",
    fill = "Types of Exposures"
  ) +
  theme_minimal() + theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



3. Name two social factors or behaviours you would expect from the sampled demographic that would increase the likelihood of exposure to MERS-CoV (particularly for the elderly population). An example could be “hospital hopping” where patients visit multiple hospitals to see different medical professionals.

The main hospitals where MERS-CoV spread were Samsung Medical Center (52.5%) and Pyeongtaek St. Mary’s Hospital (PTSM) (22.8%). These two hospitals accounted for 74% of all infections. In the case of the elderly population (49 cases), 20 of them were infected at Samsung Medical Center.

Hospital Room visits where family members and HCW could come and go could increase the likelihood of exposure to MERS-CoV. Another behavior that could increase the likelihood of becoming exposed would be to go to the emergency room for another medical condition.

```
# understanding the elderly population's behavior
elderly <- linelist %>%
  filter(age >= 65)
nrow(elderly)
```

```
## [1] 49
```

```
hospital_elderly <- elderly %>%
  group_by(loc_hosp) %>%
  summarise(
    count = n(),
```

```

percentage = count / nrow(linelist) * 100,
mean_age = mean(age, na.rm= TRUE),
.groups = 'drop'
) %>%
  arrange(desc(count))
head(hospital_elderly)

```

```

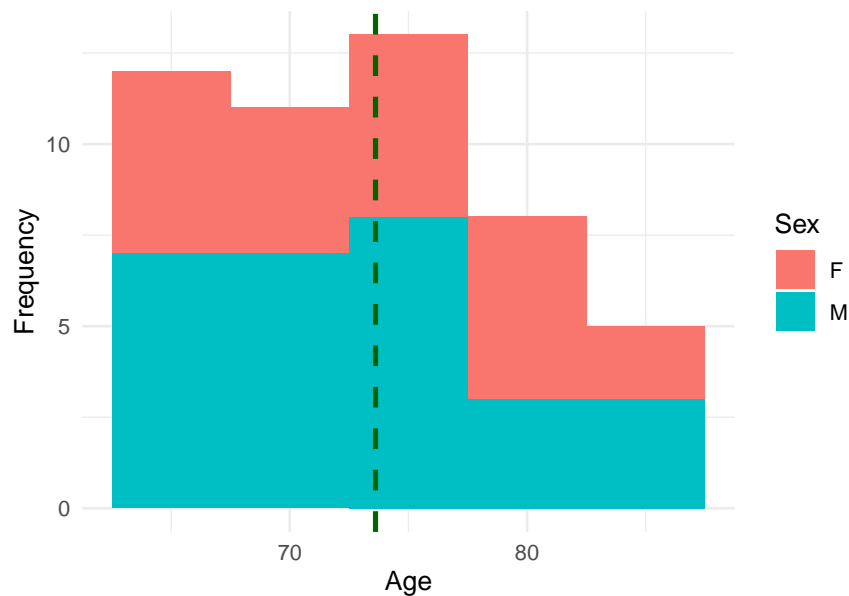
## # A tibble: 6 x 4
##   loc_hosp                count percentage mean_age
##   <chr>                  <int>      <dbl>    <dbl>
## 1 Samsung Medical Center      20      12.3     72.4
## 2 Konyang University Hospital    9       5.56     74.1
## 3 Dae Cheong Hospital          7       4.32     76.9
## 4 Pyeongtaek St. Mary          7       4.32     74.3
## 5 Pyeongtaek goodmorning hospital 4       2.47     73.2
## 6 Hallym University Medical Center 2       1.23     71

```

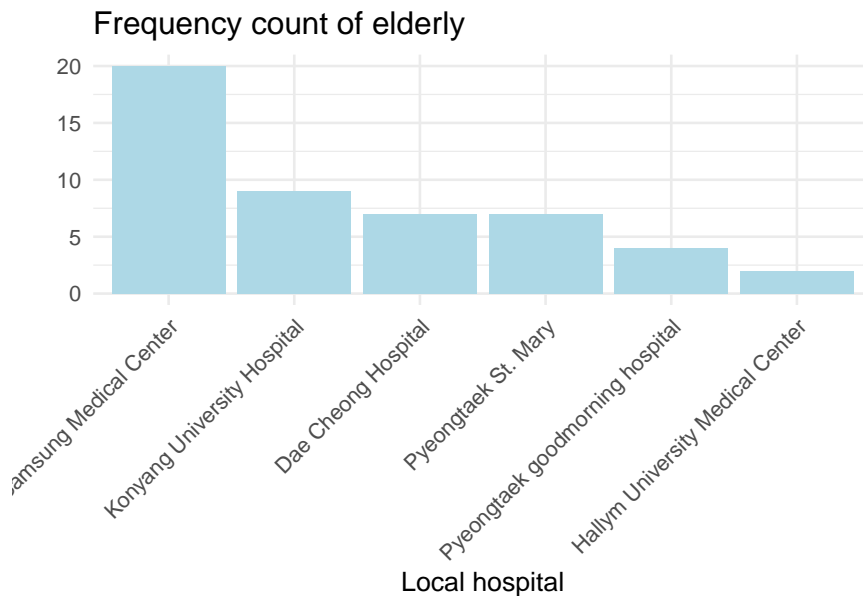
```

elderly_dist <- ggplot(elderly) +
  geom_histogram(aes(x=age, fill=sex), binwidth=5) +
  labs(
    x="Age",
    y="Frequency",
    fill="Sex"
  ) +
  geom_vline( # Average age as vertical line
    xintercept = mean(elderly$age),
    color = "darkgreen",
    linetype = "dashed",
    linewidth = 1 ) +
  theme_minimal()
elderly_dist

```



```
# distribution across hospitals (elderly)
ggplot(elderly) +
  geom_bar(
    aes(x = fct_infreq(loc_hosp) ),
    fill = "lightblue" ) +
  labs(
    title = "Frequency count of elderly",
    x = "Local hospital",
    y = ""
  ) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) # changing the text's angle
```



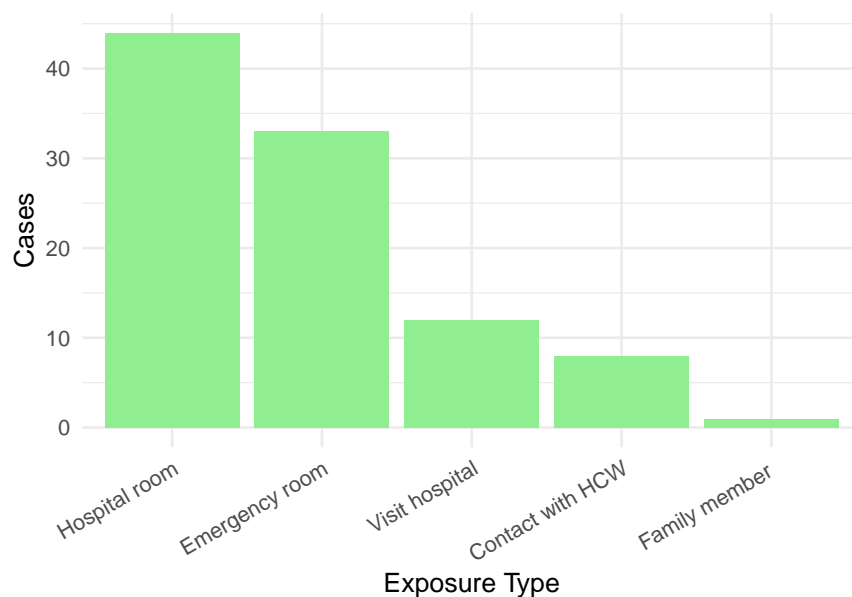
4. Synthesize your statistical finding with your identified social factors and behaviours to explain how **social status and norms** acted as an **amplifier** of the biological risk posed by age. Maybe it's a good time to look into Korean culture too!
5. Repeat task (d) but with the following points: saturated emergency rooms and family caregiving, common in Korea, where there is a cultural expectation for family members (usually younger women or older teens) to stay in the patient's room providing simple emotional and practical support.

```
exposures <- contacts %>%
  group_by(exposure) %>%
  summarise (
    count = n(),
    percentage = count / nrow(linelist) * 100 %>%
  )
  arrange(desc(count))
head(exposures)
```

```
## # A tibble: 5 x 3
##   exposure      count percentage
##   <fct>         <int>      <dbl>
```

```
## 1 "Hospital room"      44      27.2
## 2 "Emergency room"    33      20.4
## 3 "Visit hospital"    12       7.41
## 4 "Contact with HCW"   8       4.94
## 5 "Family member "    1       0.617
```

```
exposure_types <- ggplot(contacts) +
  geom_bar(aes(x = fct_infreq(exposure)), fill="lightgreen") +
  labs(
    x="Exposure Type",
    y="Cases" ) +
  theme_minimal() + theme(axis.text.x = element_text(angle = 30, hjust = 1))
exposure_types
```



2. Institutional Analysis

1. Calculate the median reporting lag from the original variables provided. On average, how many days were people sick and potentially interacting with their social networks before the institution ‘captured’ them?

The median reporting between the date of onset and the date of reporting is 5 days. For 5 days, cases were interacting with other people, potentially transmitting MERS-CoV. SK_1, one of the super-spreaders and responsible for 26 infections, had a lag of 16 days between the onset of symptoms and their report.

```
median_reporting <- linelist %>%
  mutate (
    lag_report = dt_report - dt_onset,
    lag_max = dt_report - dt_start_exp, # max delay
    lag_min = dt_report - dt_end_exp, # min delay
    lag_exp_to_onset = lag_max - lag_report, # days between exposure and reporting
    week = week(dt_report)) %>%
  filter(
```

```

!is.na(lag_report),
!is.na(lag_max),
!is.na(lag_min),
!is.na(lag_exp_to_onset)) %>%
group_by(week) %>%
summarise(
  median_max = median(lag_max, na.rm = TRUE),
  median_min = median(lag_min, na.rm = TRUE),
  median_report = median(lag_report, na.rm = TRUE),
  median_exposures = median(lag_max, na.rm = TRUE),
  .groups = 'drop')
head(median_reporting)

```

```

## # A tibble: 5 x 5
##   week median_max median_min median_report median_exposures
##   <dbl> <drtn>      <drtn>      <drtn>      <drtn>
## 1    20  5.0 days    4 days    2.0 days    5.0 days
## 2    21 10.0 days    8 days    1.5 days   10.0 days
## 3    22 16.0 days   14 days    8.0 days   16.0 days
## 4    23 12.5 days   10 days    5.0 days   12.5 days
## 5    24 16.0 days   15 days    4.0 days   16.0 days

```

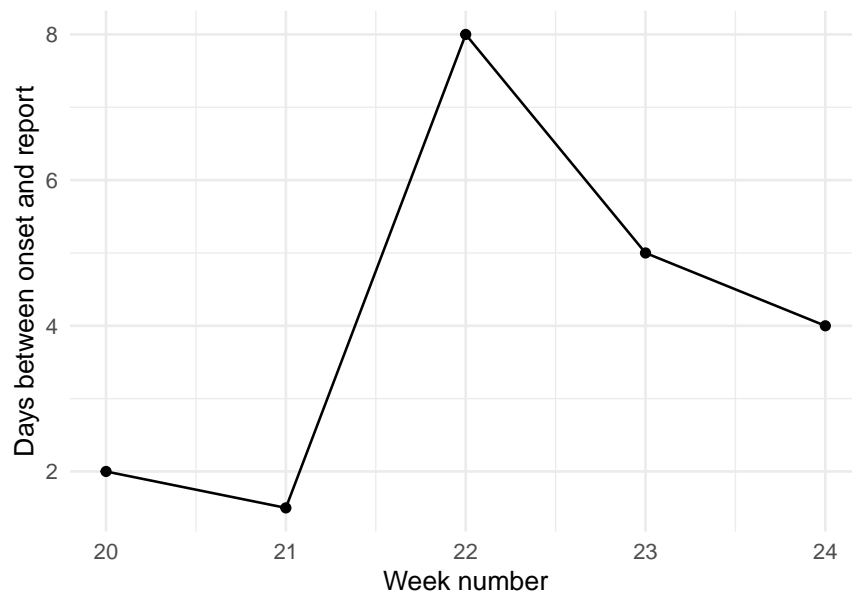
2. Analyse your new variable as a function of time. What does the trend tell you about the institution's ability to 'learn' and mobilize resources over time?

In week 22, the median reporting lag between the onset of symptoms and the report was the highest recorded (8 days). After week 22, the Korea CDC learned and the cases reported in weeks 23 and 24 decreased their lags to 5 and 4 days respectively.

```

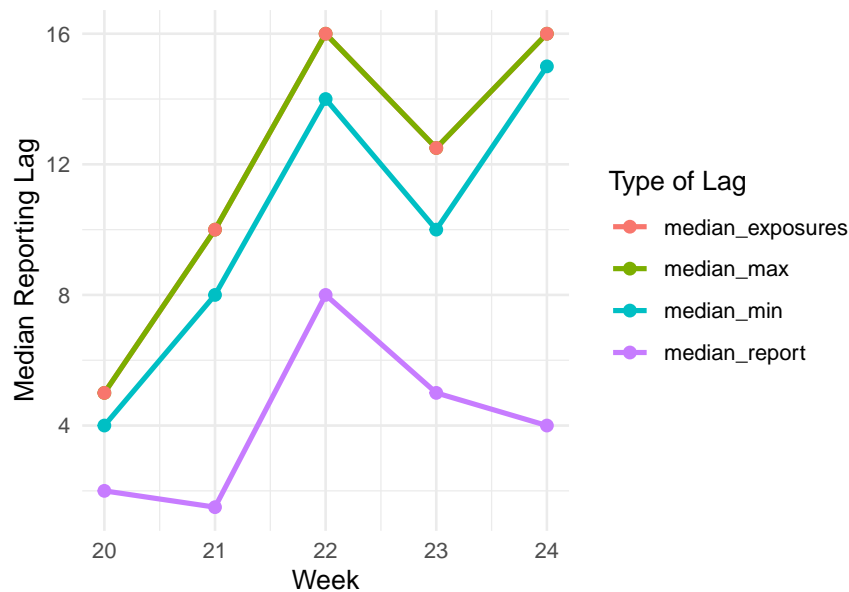
# median lag across time
ggplot(median_reporting, aes(x=week, y=median_report)) +
  geom_line() + geom_point() + theme_minimal() +
  labs(x="Week number", y="Days between onset and report")

```



```
# analyzing different lags
lag_by_week_long <- median_reporting %>%
  pivot_longer(
    cols = starts_with("median"),
    names_to = "lag_type",
    values_to = "values" )

ggplot(lag_by_week_long,
  aes(x = week,
      y = values,
      color = lag_type,
      group = lag_type) ) +
  geom_line(linewidth = 1) +
  geom_point(size = 2) +
  labs(x="Week", y="Median Reporting Lag", color="Type of Lag")+ theme_minimal()
```



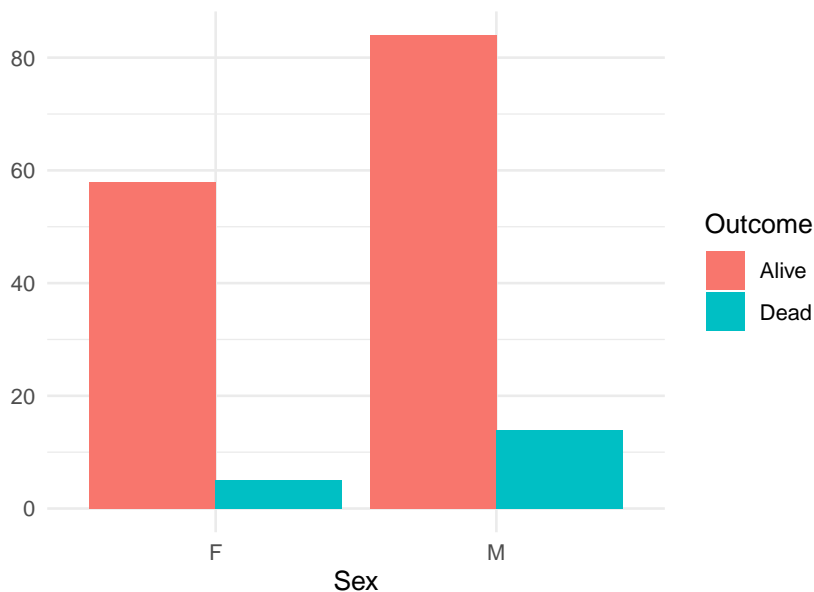
3. In MERS-CoV outbreaks in the Middle East, men were infected much more often. In this South Korean outbreak, is there a gender imbalance? What could this imply about the gender demographics of caregiving in Korean hospitals?

In this outbreak, men were infected 1.5 times higher than women.

```
# looking up the gender distribution
gender_stats <- linelist %>%
  filter(
    place_infect != "Middle East"
  ) %>%
  group_by(place_infect, sex, outcome) %>%
  summarise(
    count = n(),
    percentage = count / nrow(linelist) * 100,
    median_age = median(age, na.rm = TRUE))
gender_stats
```

```
## # A tibble: 4 x 6
## # Groups:   place_infect, sex [2]
##   place_infect      sex outcome count percentage median_age
##   <fct>          <fct> <fct>   <int>      <dbl>      <dbl>
## 1 Outside Middle East F    Alive     58      35.8       55.5
## 2 Outside Middle East F    Dead       5       3.09       68
## 3 Outside Middle East M    Alive    84      51.9       54.5
## 4 Outside Middle East M    Dead     14       8.64       68
```

```
ggplot(gender_stats, aes(x= fct_infreq(sex), y=count, fill=outcome)) +
  geom_bar(stat="identity", position="dodge") + theme_minimal() +
  labs(x="Sex", y="", fill="Outcome")
```



Bonus - investigate why there were more male than female cases in the middle east spread of MERS-CoV. Why would this not apply in Korea?

3. The Super-Spreader

- Using the appropriate dataset, calculate the number of infections attributed to each case. Calculate the mean and the maximum of this variable. What is the ratio between the maximum “infector” and the average?

For this analysis, I will categorize the super-spreader as those ids that infected 5 or more other cases.

Based on the Contacts data set and the previous assumption, we can identify 10 spreaders in this outbreak, and 3 super-spreaders: SK_14, SK_1, and SK_16. These 3 super-spreaders are males with an average age of 35 years old and located in the Pyeongtaek St. Mary Hospital (PTSM).

The 3 super-spreaders were responsible for 85 infections (52%). The super-spreader infected 4 times more than the average spreaders.

```
# count of infections caused by each spreader
all_cases <- nrow(linelist)
```



```

spreader_count <- contacts %>%
  group_by(from) %>%
  summarise(
    infections = n(),
    percentage_of_infections = infections / all_cases * 100,
    .groups = 'drop' ) %>%
  arrange(desc(infections))
head(spreader_count)

```

```

## # A tibble: 6 x 3
##   from infections percentage_of_infections
##   <chr>      <int>                <dbl>
## 1 SK_14      38                23.5
## 2 SK_1       26                16.0
## 3 SK_16      21                13.0
## 4 SK_15       4                 2.47
## 5 SK_6        2                 1.23
## 6 SK_76       2                 1.23

```

```

# how did these spreaders infect the other cases?
superspreaders <- contacts %>%
  group_by(from,
            exposure) %>%
  summarise(
    infections = n(),
    percentage_of_infections = infections / all_cases * 100,
    .groups = 'drop') %>%
  arrange(desc(infections))
head(superspreaders)

```

```

## # A tibble: 6 x 4
##   from exposure      infections percentage_of_infections
##   <chr> <fct>          <int>                <dbl>
## 1 SK_14 Emergency room      30                18.5
## 2 SK_16 Hospital room     20                12.3
## 3 SK_1  Visit hospital     12                 7.41
## 4 SK_1  Hospital room       9                 5.56
## 5 SK_14 Hospital room       6                 3.70
## 6 SK_1  Contact with HCW     4                 2.47

```

```

# What is the ratio between the maximum "infector" and the average?
# formula for a ratio = max / average

```

```

# calculating the max
max_infector_count <- max(spreader_count$infections, na.rm = TRUE)
max_infector_count

```

```
## [1] 38
```

```
# calculating the average
avg_spreader_count <- mean(spreader_count$infections, na.rm = TRUE)
avg_spreader_count
```

```
## [1] 8.909091
```

```
# calculating the ratio
ratio <- max_infector_count / avg_spreader_count
ratio
```

```
## [1] 4.265306
```

- b) Visualise the distribution of secondary cases. What does the shape of your plot imply about how “democratic” or “equal” disease transmission was in this specific sample?
- c) Summarise the demographic profile of the “super spreader” using both datasets. Based on where they were infected and their demographic profile, hypothesize a social reason why this specific individual had such a high number of secondary cases.

```
# left join to dig deeper into the superspreaders' details
superspreaders_deep <- superspreaders %>%
  rename(id = from) %>% # rename 'from' to 'id'
  left_join(linelist, by = "id") %>%
  select(id,
         age,
         age_class,
         sex,
         infections,
         percentage_of_infections,
         place_infect,
         reporting_ctry,
         exposure,
         loc_hosp,
         outcome) %>%
  filter(infections >= 5)
superspreaders_deep
```

```
## # A tibble: 5 x 11
##   id      age age_class sex   infections percentage_of_infections place_infect
##   <chr> <int> <chr>    <fct>      <int>          <dbl> <fct>
## 1 SK_14   35 30-39     M           30          18.5 Outside Middl~
## 2 SK_16   40 40-49     M           20          12.3 Outside Middl~
## 3 SK_1    68 60-69     M           12           7.41 Middle East
## 4 SK_1    68 60-69     M           9            5.56 Middle East
## 5 SK_14   35 30-39     M           6            3.70 Outside Middl~
## # i 4 more variables: reporting_ctry <fct>, exposure <fct>, loc_hosp <chr>,
## #   outcome <fct>
```

Bonus - Test the “Pareto Principle” (the 20/80 rule) on this data. Calculate what percentage of the total infections were caused by the top 20% of infectors. Does this outbreak fit the rule?

Appendix

Infection Network

```
library(igraph)
links <- data.frame(
  source= contacts$from,
  target = contacts$to
)
network <- graph_from_data_frame(
  d=links,
  directed=F
)

plot(network,
  vertex.size = 10,
  vertex.label.cex = 0.5,
  vertex.color="lightblue")
```

