```
Midterm 2 W24
Jacques Mak
2024-02-27
Instructions
Answer the following questions and complete the exercises in RMarkdown. Please embed all of your code and push your final work to your
repository. Your code must be organized, clean, and run free from errors. Remember, you must remove the # for any included code chunks to
run. Be sure to add your name to the author header above.
Your code must knit in order to be considered. If you are stuck and cannot answer a question, then comment out your code and knit the
document. You may use your notes, labs, and homework to help you complete this exam. Do not use any other resources- including Al
assistance.
Don't forget to answer any questions that are asked in the prompt. Some questions will require a plot, but others do not-make sure to read each
question carefully.
For the questions that require a plot, make sure to have clearly labeled axes and a title. Keep your plots clean and professional-looking, but you
are free to add color and other aesthetics.
Be sure to follow the directions and upload your exam on Gradescope.
Background
In the data folder, you will find data about shark incidents in California between 1950-2022. The data are from: State of California-Shark
Incident Database.
Load the libraries
 library("tidyverse")
 library("janitor")
 library("naniar")
Load the data
Run the following code chunk to import the data.
 sharks <- read csv("data/SharkIncidents 1950 2022 220302.csv") %>% clean names()
Questions
   1. (1 point) Start by doing some data exploration using your preferred function(s). What is the structure of the data? Where are the missing
     values and how are they represented?
 glimpse(sharks)
 ## Rows: 211
 ## Columns: 16
 ## $ incident num
                         <chr> "1", "2", "3", "4", "5", "6", "7", "8", "9", "10", "1...
 ## $ month
                         <dbl> 10, 5, 12, 2, 8, 4, 10, 5, 6, 7, 10, 11, 4, 5, 5, 8, ...
 ## $ day
                         <dbl> 8, 27, 7, 6, 14, 28, 12, 7, 14, 28, 4, 10, 24, 19, 21...
                         <dbl> 1950, 1952, 1952, 1955, 1956, 1957, 1958, 1959, 1959,...
 ## $ year
                         <chr> "12:00", "14:00", "14:00", "12:00", "16:30", "13:30",...
 ## $ time
                         <chr> "San Diego", "San Diego", "Monterey", "Monterey", "Sa...
 ## $ county
 ## $ location
                         <chr> "Imperial Beach", "Imperial Beach", "Lovers Point", "...
                         <chr> "Swimming", "Swimming", "Swimming", "Freediving", "Sw...
 ## $ mode
 ## $ injury
                         <chr> "major", "minor", "fatal", "minor", "major", "fatal",...
                         <chr> "surface", "surface", "surface", "surface"...
 ## $ depth
 ## $ species
                         <chr> "White", "White", "White", "White", "White", "White", "
 ## $ comment
                         <chr> "Body Surfing, bit multiple times on leg, thigh and b...
 ## $ longitude
                         <chr> "-117.1466667", "-117.2466667", "-122.05", "-122.15",...
 ## $ latitude
                         <dbl> 32.58833, 32.58833, 36.62667, 36.62667, 35.13833, 35....
 ## $ confirmed source <chr> "Miller/Collier, Coronado Paper, Oceanside Paper", "G...
 anyNA(sharks)
 ## [1] TRUE
 miss var summary(sharks)
 ## # A tibble: 16 × 3
       variable
                          n_miss pct_miss
        <chr>
                           <int>
                                     <dbl>
 ## 1 wfl case number
                           202 95.7
 ## 2 time
                              7
                                   3.32
 ## 3 latitude
                                    2.84
 ## 4 longitude
                                5 2.37
 ## 5 confirmed source
                               1 0.474
 ## 6 incident num
                                     0
 ## 7 month
                                    0
 ## 8 day
 ## 9 year
 ## 10 county
 ## 11 location
 ## 12 mode
 ## 13 injury
 ## 14 depth
 ## 15 species
 ## 16 comment
   2. (1 point) Notice that there are some incidents identified as "NOT COUNTED". These should be removed from the data because they were
      either not sharks, unverified, or were provoked. It's OK to replace the sharks object.
 sharks%>%
   filter(incident_num!="NOT COUNTED")
 ## # A tibble: 202 × 16
        incident_num month day year time
                                                  county
                                                              location mode injury depth
                      <dbl> <dbl> <dbl> <chr>
        <chr>
                                                  <chr>
                                                              <chr>
                                                                        <chr> <chr> <chr>
                       10
                                                  San Diego Imperia... Swim... major surf...
 ## 1 1
                                 8 1950 12:00
                     5 27 1952 14:00
 ## 2 2
                                                  San Diego Imperia... Swim... minor surf...
           12 7 1952 14:00
2 6 1955 12:00
 ## 3 3
                                                  Monterey Lovers ... Swim... fatal surf...
                                                  Monterey Pacific... Free... minor surf...
 ## 4 4
           8 14 1956 16:30
4 28 1957 13:30
 ## 5 5
                                                  San Luis ... Pismo B... Swim... major surf...
 ## 6 6
                                                  San Luis ... Morro B... Swim... fatal surf...
                   10 12 1958 Unknown San Diego Coronad... Swim... major surf...
5 7 1959 17:30 San Franc... Baker B... Swim... fatal surf...
 ## 7 7
 ## 8 8
                       6 14 1959 17:00
 ## 9 9
                                                  San Diego La Jolla Free... fatal surf...
 ## 10 10
                                28 1959 19:30
                                                  San Diego La Jolla Free... minor surf...
 ## # i 192 more rows
 ## # i 6 more variables: species <chr>, comment <chr>, longitude <chr>,
 ## # latitude <dbl>, confirmed_source <chr>, wfl_case_number <chr>
   3. (3 points) Are there any "hotspots" for shark incidents in California? Make a plot that shows the total number of incidents per county.
     Which county has the highest number of incidents?
 sharks %>%
   filter(incident_num!="NOT COUNTED") %>%
   count(county)%>%
   arrange(desc(n))
 ## # A tibble: 21 × 2
        county
                             n
        <chr>
                         <int>
 ## 1 San Diego
    2 Santa Barbara
 ## 3 Humboldt
                            18
    4 San Mateo
                          18
 ## 5 Marin
 ## 6 Monterey
                      15
 ## 7 Santa Cruz
 ## 8 Sonoma
 ## 9 San Luis Obispo
 ## 10 Los Angeles
 ## # i 11 more rows
##San Diego has the largest total number of incident.
 sharks %>%
   ggplot(aes(fill = county, x = county)) +
   geom_bar() +
   coord_flip() +
   labs(title = "Number of Incidents per County",
         y = "Number of Incidents",
         x = "County") +theme_minimal()
                     Number of Incidents per County
              Ventura
              Sonoma
            Santa Cruz
         Santa Barbara
                                                      county
            San Mateo
                                                          Del Norte
                                                                                Mendocino
        San Luis Obispo
         San Francisco
                                                          Humboldt
                                                                                Monterey
            San Diego
                                                          Island - Catalina
                                                                                Orange
              Orange
                                                          Island - Farallones
                                                                                San Diego
             Monterey
                                                          Island - San Miguel
                                                                                San Francisco
            Mendocino
                Marin
                                                          Island - San Nicolas
                                                                                San Luis Obispo
           Los Angeles
                                                          Island - Santa Barbara
                                                                                San Mateo
     Island - Santa Rosa
                                                          Island - Santa Cruz
                                                                                Santa Barbara
      Island - Santa Cruz
                                                          Island - Santa Rosa
                                                                                Santa Cruz
   Island - Santa Barbara
     Island - San Nicolas
                                                          Los Angeles
                                                                                Sonoma
      Island - San Miguel
                                                          Marin
                                                                                Ventura
      Island - Farallones
        Island - Catalina
             Humboldt
             Del Norte
                                10
                                     15
                          Number of Incidents
   4. (3 points) Are there months of the year when incidents are more likely to occur? Make a plot that shows the total number of incidents by
      month. Which month has the highest number of incidents?
 sharks$month <- as.factor(sharks$month)</pre>
 sharks %>%
   ggplot(aes(fill = month, x = month)) +
   geom_bar() +
   coord_flip() +
   labs(title = "Shark Incidents by Month",
         y = "Number of Incidents",
        x = "County") +theme_minimal()
      Shark Incidents by Month
   12
   11
                                                                                      month
   10
    8
 County
                                                                                              ##October looks like a hotspot of the
    6
    5
                                                                                          11
                                                                                          12
                            10
                                                                   30
                                   Number of Incidents
incidents.
   5. (3 points) How do the number and types of injuries compare by county? Make a table (not a plot) that shows the number of injury types by
     county. Which county has the highest number of fatalities?
    sharks %>%
   group_by(county, injury) %>%
    summarise(total = n()) %>%
   pivot_wider(names_from = injury,
                 values_from = total) %>%
   arrange(desc(fatal))
 ## `summarise()` has grouped output by 'county'. You can override using the
 ## `.groups` argument.
 ## # A tibble: 22 × 8
 ## # Groups: county [22]
        county
                              minor none major `minor*` `none*` fatal `major*`
                              <int> <int> <int>
                                                             <int> <int>
        <chr>
                                                     <int>
                                                                              <int>
     1 San Luis Obispo
                                        7
                                            3
                                                        NA
                                                                NA
                                                                        3
     3 San Diego
                                         9
                                                         1
                                                                 1
                                                                                 NA
     4 Santa Barbara
                                        9
                                                        NA
                                                                NA
                                                                                 NA
      5 Island - San Miguel
                                       NA
                                               2
                                                        NA
                                                                NA
                                                                        1
                                                                                 NA
      6 Los Angeles
                                              NA
                                                        NA
                                                                NA
                                                                                  1
     7 Mendocino
                                       NA
                                               3
                                                        NA
                                                                NA
                                                                        1
                                                                                 NA
      8 San Francisco
                                 NA
                                        1
                                              NA
                                                        NA
                                                                NA
                                                                                 NA
     9 San Mateo
                                       12
                                                        1
                                                                NA
                                                                                 NA
 ## 10 Santa Cruz
                                                        NA
                                                                NA
                                                                                 NA
 ## # i 12 more rows
##San luis Obispo has to most fatal incidents.
   6. (2 points) In the data, mode refers to a type of activity. Which activity is associated with the highest number of incidents?
 sharks %>%
   ggplot(aes(fill = mode, x = mode)) +
   geom_bar() +
   coord_flip() +
   labs(title = "Number of Incidents per Mode",
         y = "Number of Incidents",
         x = "County")+theme_minimal()
                   Number of Incidents per Mode
     Walking in shallow
           Swimming
                                                                        mode
     Surfing / Boarding
                                                                             Freediving
                                                                             Hookah Diving
         Scuba Diving
                                                                             Kayaking / Canoeing
 County
                                                                             Paddleboarding
                                                                                              ##Surfing/ Boardinf is associating
                                                                             Scuba Diving
       Paddleboarding
                                                                             Surfing / Boarding
                                                                             Swimming
   Kayaking / Canoeing
                                                                             Walking in shallow
        Hookah Diving
           Freediving
                               20
                                                                80
                                   Number of Incidents
with the most incidents.
   7. (4 points) Use faceting to make a plot that compares the number and types of injuries by activity. (hint: the x axes should be the type of
     injury)
 sharks %>%
   group_by(mode, injury) %>%
   ggplot(aes(x = injury)) +
   geom_bar(fill = "light green") +
   facet_wrap(~mode) +
   labs(title = "Types of Shark Injuries by Mode",
         x = "Injury Type",
         y = "Number of Incidents") +
   theme_minimal()+ theme(axis.text.x = element_text(angle = 30, hjust=1))
      Types of Shark Injuries by Mode
                Freediving
                                            Hookah Diving
                                                                       Kayaking / Canoeing
   30
   20
   10
    0
                                                                        Surfing / Boarding
              Paddleboarding
                                             Scuba Diving
 Number of Incidents
   10
                                                                fatal major minor minor none none
                 Swimming
                                           Walking in shallow
   30
   20
   10
         major minor minor none none
                                            Injury Type
   8. (1 point) Which shark species is involved in the highest number of incidents?
 sharks %>%
   ggplot(aes(fill = species, x = species)) +
   geom_bar() +
   coord_flip() +
   labs(title = "Number of Incidents by species",
         y = "Number of Incidents",
         x = "County") +theme_minimal()
              Number of Incidents by species
         White
      Unknown
                                                                              species
       Thresher
                                                                                  blue
                                                                                  Blue
       Sevengill
                                                                                  Blue*
        Salmon
                                                                                  Hammerhead
                                                                                  Killer Whale
 County
         Mako
                                                                                  Leopard
                                                                                              The white shark has contributes the
       Leopard
                                                                                  Mako
                                                                                  Salmon
     Killer Whale
                                                                                  Sevengill
   Hammerhead
                                                                                  Thresher
                                                                                  Unknown
         Blue*
                                                                                  White
          Blue
          blue
                                                             150
                               50
                                              100
                                   Number of Incidents
most incidents.
   9. (3 points) Are all incidents involving Great White's fatal? Make a plot that shows the number and types of injuries for Great White's only.
 sharks %>%
  filter(species=="White")%>%
   group_by(mode, injury) %>%
   ggplot(aes(x = injury)) +
   geom bar(fill = "light green") +
   facet_wrap(~mode) +
   labs(title = "Great White Sharks injury",
         x = "Injury Type",
         y = "Number of Incidents") +
   theme_minimal()+ theme(axis.text.x = element_text(angle = 30, hjust=1))
      Great White Sharks injury
                Freediving
                                            Hookah Diving
                                                                       Kayaking / Canoeing
   30
   20
   10
    0
               Paddleboarding
                                             Scuba Diving
                                                                        Surfing / Boarding
 Number of Incidents
   30
   20
                                                                                              ##Not all of the Great White Sharks
                                                                     major minor
                 Swimming
   30
   20
   10
                                            Injury Type
incidents are fatal.
Background
Let's learn a little bit more about Great White sharks by looking at a small dataset that tracked 20 Great White's in the Fallaron Islands. The data
are from: Weng et al. (2007) Migration and habitat of white sharks (Carcharodon carcharias) in the eastern Pacific Ocean.
Load the data
 white_sharks <- read_csv("data/White sharks tracked from Southeast Farallon Island, CA, USA, 1999 2004.csv", na =
 c("?", "n/a")) %>% clean names()
 10. (1 point) Start by doing some data exploration using your preferred function(s). What is the structure of the data? Where are the missing
      values and how are they represented?
 summary(white_sharks)
                          tagging_date
 ##
         shark
                                               total_length_cm
                                                                     sex
                          Length:20
      Length:20
                                               Min.
                                                     :360.0
                                                                Length:20
                          Class :character
      Class :character
                                               1st Qu.:400.5
                                                                Class :character
                          Mode :character
      Mode :character
                                               Median:434.5
                                                                Mode :character
 ##
                                                     :436.1
                                               Mean
 ##
                                               3rd Qu.:457.0
 ##
                                                       :530.0
                                               Max.
 ##
        maturity
                          pop_up_date
                                                                   longitude
                                                 track_days
 ## Length:20
                          Length:20
                                               Min. : 14.0 Min. :-156.8
     Class: character Class: character 1st Qu.: 85.0 1st Qu.:-137.8
      Mode :character Mode :character Median :182.0
                                                                Median :-133.2
 ##
                                               Mean :166.8 Mean :-120.3
                                                                3rd Qu.:-124.3
 ##
                                               3rd Qu.:216.8
 ##
                                               Max. :367.0 Max. : 131.7
                                                                NA's :1
         latitude
                         comment
 ## Min. :20.67 Length:20
     1st Qu.:22.48
                       Class :character
      Median :26.39
                       Mode :character
      Mean :28.24
      3rd Qu.:36.00
           :38.95
     Max.
 ## NA's :1
 miss var summary(white sharks)
 ## # A tibble: 10 × 3
        variable n_miss pct_miss
        <chr> <int>
                                    <dbl>
     1 sex 3
2 maturity 1
    1 sex
                                       15
 ## 3 longitude
 ## 4 latitude
 ## 5 shark 0
## 6 tagging_date 0
 ## 7 total_length_cm 0
## 8 pop_up_date 0
 ## 8 pop_up_date
 ## 9 track_days
 ## 10 comment
 11. (3 points) How do male and female sharks compare in terms of total length? Are males or females larger on average? Do a quick search
     online to verify your findings. (hint: this is a table, not a plot).
 white_sharks %>%
   filter(sex!="NA") %>%
   group_by(sex) %>%
    summarize(mean_lenth=mean(total_length_cm, na.rm=T))
 ## # A tibble: 2 × 2
       sex mean_lenth
       <chr>
                   <dbl>
                    462
 ## 1 F
 ## 2 M
                    425.
 12. (3 points) Make a plot that compares the range of total length by sex.
 white_sharks %>%
   filter(sex!="NA") %>%
   ggplot(aes(x=sex, y=total_length_cm, fill=sex)) +
   geom_boxplot(na.rm=T, alpha=0.5) +
   labs(title="Range of Length by sex",
         x=NULL,
         y="Lenth",
         fill="sex") +
   theme(plot.title = element_text(size=12, face="bold"),
          axis.title.x = element_text(size=10),
          axis.title.y = element text(size=10))+theme minimal()
       Range of Length by sex
   500
                                                                                      sex
                                                                                      ⊨ F
                                                                                      і м
   400
 13. (2 points) Using the sharks or the white_sharks data, what is one question that you are interested in exploring? Write the question and
      answer it using a plot or table.
What are the number and types of injuries specifically in Sonoma?
 sharks %>%
  filter(county=="Sonoma")%>%
   group_by(mode, injury) %>%
   ggplot(aes(x = injury)) +
    geom_bar(fill = "light green") +
    facet_wrap(~mode) +
   labs(title = "Sonoma injury",
         x = "Injury Type",
         y = "Number of Incidents") +
   theme_minimal()+ theme(axis.text.x = element_text(angle = 30, hjust=1))
     Sonoma injury
                Freediving
                                          Kayaking / Canoeing
                                                                          Scuba Diving
 Number of Incidents
```

Surfing / Boarding

minor

2

Swimming

Injury Type