

STA303/1002: Mixed assessment 1

Starship crew analysis

Chief Science Officer [your name]; ID: [your student ID]

Information	Note
Name	Mixed assignment 1
Type	Type 2
Value	14%
Due	This untimed submission must be submitted by 4:30 p.m. ET Wednesday, Feb 24
Submission	PDF & RMD: https://q.utoronto.ca/courses/204826/assignments/415116
link	
Accommodations and extension policy	If you miss a type 2 assessment due to illness or a serious personal emergency, please complete this form within ONE week of the due date of the assignment (i.e. the end of the timed assessment window).

Mixed assessment 1 has two components:

- Untimed guided analysis (this)
- [Timed assessment](#) (50 minutes; 24-hour assessment window is 4:30 p.m. ET Tuesday, Feb 23 to 4:30 p.m. ET Wednesday, Feb 24)
- See the [mixed assessments overview page](#) for further information and revisions links.

Instructions

1. Update the `yml` at the top of this document to have your name and your student ID. There are TWO places you need to do this for each one, probably on lines 4 and 12. I.e., replace the square brackets and everything inside them with the appropriate details. Your student ID is all numbers (usually 10, sometimes 8 or 9), it is NOT your UTORid.
2. Complete the guided analysis below. You will want to complete this BEFORE attempting your timed assessment.
3. Complete your [timed assessment](#). It will require your work in this document, as well general STA303 content knowledge.
4. Once you've written your code and are ready to knit change `knitr::opts_chunk$set(eval = FALSE)` to read `eval = TRUE` in the libraries chunk.
5. Knit this .Rmd to .pdf and submit BOTH files to the submission link in the table above.

Note: This component is ungraded BUT there is a 10-percentage point penalty for not submitting your .Rmd and .pdf to the Quercus dropbox by the end of the assessment window, i.e., by Wednesday, February 24 at 4:30 p.m. ET. The intention is to allow confirmation of your work for academic integrity purposes and/or if as a way to confirm your personalized data if there are issues.

Setting up your libraries

If you are working on this on the Jupyter Hub, the `tidyverse`, `devtools`, `lme4`, `lattice` and `lmtest` packages will already be installed. If you're working locally, you'll have to install them first if they are not already installed. You will also need to install the `randomNames` package and the `myStarship` package from GitHub. All the code you need to do this is in the `setup` chunk below.

```
# RUN THIS CHUNK FIRST! You should only need to run it once on your local machine.
# On the JupyterHub, you may need to run it at the beginning of each new session.

# These are the packages you will need for this activity.
packages_needed <- c("tidyverse", "devtools", "lme4",
                    "lattice", "lmtest", "randomNames")

package.check <- lapply(
  packages_needed,
  FUN = function(x) {
    if (!require(x, character.only = TRUE)) {
      install.packages(x, dependencies = TRUE,
        repos = "https://cloud.r-project.org/") # you may need to change the mirror if
        # you're in China (and potentially other countries.)
        # Students in China have reported that
        # "https://mirrors.tuna.tsinghua.edu.cn/CRAN/" worked for them.
    }
  }
)

# Remove objects no longer needed
rm(packages_needed, package.check)

# You may be prompted to install or update additional packages
# If so, you'll see a message in the console
# Type a enter/return in the console to skip updating
devtools::install_github("elb0/myStarship", force = TRUE)

# Run libraries for easy access to the functions we'll be using
library(tidyverse)
library(lme4)
library(myStarship)

# Once you've updated the code and are ready to knit
# change this to eval = TRUE
knitr::opts_chunk$set(eval = FALSE)
```

Get your data

IMPORTANT you must update your student ID in the function in the following chunk. You will be graded based on your *unique dataset* and so risk losing extensive marks if you use the wrong dataset.

```
# put your student ID in here
get_my_starship(111111111)

# after you run this function, your unique dataset will appear in the environment
# it will be called crew_data
```

The goal

You are the Chief Science Officer of the SS Gres. You have data about the productivity of the crew over a 12 week period after a shore leave (a holiday break for the crew). For each member of the crew you also have data on their **rank** within Starfleet, their role on the ship (**position**), which of the three main divisions (**division**) they are in (Command, Operations, Science), as well as their sub-division (**sub_division**, e.g. Engineering is a sub-division of Operations). You also know their **gender** (Feminine, Masculine, Non-binary), **name**, what their GPA upon graduating from Starfleet Academy was (**starfleet_gpa**, 0-10 scale, 10 being the best grade), their perseverance score (**perseverance_score**) from their most recent psych assessment (0-10 scale, 10 being high perseverance). **week** indicates the weeks since the shore leave (1 to 12) and their **productivity** score for each week is recorded.

Each crewmember is assigned to a duty shift (**duty_shift**). There are four 8-hour shifts covering each 24 hour period, Alpha, Beta, Delta and Gamma. Within each duty shift, each crewmember is assigned to a team (**shift_team**). Teams are numbered 1 to 6, or sometimes fewer, and these labels aren't meaningful, they are just for administrative purposes. E.g., being Team 1 in Alpha shift has nothing to do with being Team 1 in Beta shift.

The crewmembers in Team 2 on the Gamma shift are assigned to work together as a unit, but they are only considered to be 'working' with other members of Team 2 on Gamma shift, not the rest of the Gamma shift, nor the crew in Team 2 of other shifts.

Your goal is to better understand productivity aboard your ship.

```
glimpse(crew_data)
```

```
## Rows: 3,012
## Columns: 13
## $ crew_id      <dbl> 42229, 42229, 42229, 42229, 42229, 42229, 42229, ...
## $ rank         <chr> "Captain", "Captain", "Captain", "Captain", "Cap...
## $ position     <chr> "Captain", "Captain", "Captain", "Captain", "Cap...
## $ division     <chr> "Command", "Command", "Command", "Command", "Com...
## $ sub_division <chr> "Command", "Command", "Command", "Command", "Com...
## $ gender       <chr> "Feminine", "Feminine", "Feminine", "Feminine", ...
## $ name         <chr> "Sydni Gates", "Sydni Gates", "Sydni Gates", "Sy...
## $ duty_shift   <chr> "Alpha", "Alpha", "Alpha", "Alpha", "Alpha", "Al...
## $ shift_team   <chr> "Team 1", "Team 1", "Team 1", "Team 1", "Team 1"...
## $ starfleet_gpa <dbl> 7.82, 7.82, 7.82, 7.82, 7.82, 7.82, 7.82, 7.82, ...
## $ perseverance_score <dbl> 7.24, 7.24, 7.24, 7.24, 7.24, 7.24, 7.24, 7.24, ...
## $ week         <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 1, 2, 3, ...
## $ productivity <dbl> 41.75715, 42.63031, 42.82516, 42.11529, 43.35366...
```

Task set 1: familiarize yourself with the data

1. What is the name of your ship? Hint: check out the object `ship_name`.
2. What is the name of the Communications Officer?
3. How many crewmembers are in this dataset?

Task set 2: create/alter variables

1. The Records Officer lets you know that there is a typo in the crew dataset, where 'Engineering' has been misspelled somewhere, (maybe in one of the position titles?) but unfortunately they can't remember where or how. Find the mistake, fix it (and save that fix in the original `crew_data`) and then calculate what proportion of people in the Engineering subdivision have 'engineer' or 'engineering' in their position title.
2. Create a new variable in `crew_data` called `full_team` that indicates both the duty shift and the team each person is assigned to.
 - You may find the `str_c()` function useful.
 - You can specify how the values you're sticking together are separated with the `sep` parameter, e.g., `str_c(var1, var2, sep = " ")` would put a space between the values of `var1` and `var2` when sticking them together.
 - Don't forget that `mutate()` helps you make new variables.

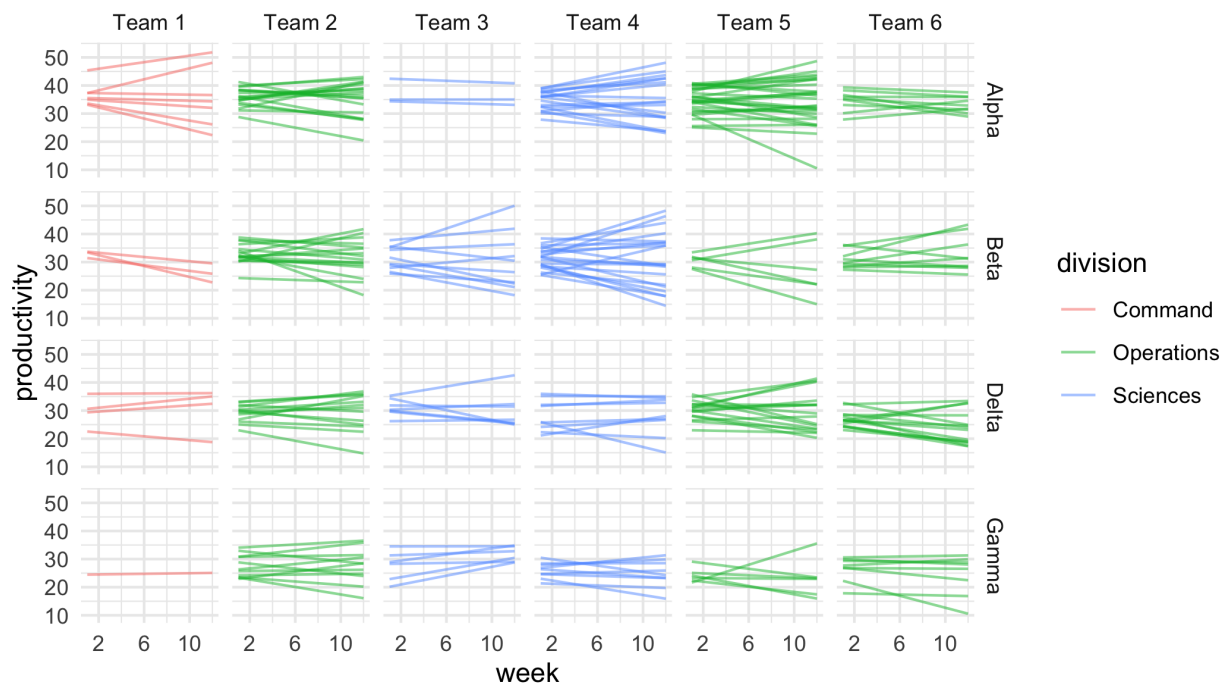
Task set 3: exploring week 1 data

1. Create a new dataset called `week1` that filters to only the observations for week 1. You must also reverse the levels of the `duty_shift` factor in `week1` so that the order is: Gamma, Delta, Beta, Alpha. You can test if you've achieved this by running `table(week1$duty_shift)`. The table should be ordered with Gamma first.
2. Using the `week1` dataset you created, create a plot with `productivity` on the y-axis, `duty_shift` on the x-axis and coloured boxplots for each `shift_team`. Use the "Dark2" colour palette from colour brewer.
 - `geom_boxplot()` is the geometry that creates boxplots.
 - use the colour aesthetic to get different boxplots for each `shift_team`
 - `scale_colour_brewer()` will allow you to choose the Dark2 palette (when completed appropriately).
3. Using the `week1` data, fit a linear model called `w1_shift` where `productivity` is the response and `duty_shift` is the only predictor. Run `summary` and `confint` on the model.
4. Fit three additional linear models and run summaries on them:
 - Name the first model `w1_team`. It should have `productivity` as the response and then `shift_team` as the only predictor. `week1` is still the data to use.
 - Name the first model `w1_int`. It should have `productivity` as the response and then the main effects and interaction of `duty_shift` and `shift_team` as the predictors. `week1` is still the data to use.
 - Name the second model `w1_full`. It should have `productivity` as the response and `full_team` as the only predictor. `week1` is still the data to use.

Task set 4: productivity post shore leave

1. Replace the 1s and add whatever other aesthetics are required in the `ggplot()` function to recreate the graph below for your particular ship. Note that each line represents the productivity trend for one crewmember over the 12 week period.

```
crew_data %>%
  ggplot(aes(y = 1, x = 1)) +
  geom_line(stat="smooth", method = "lm", formula = 'y~x', alpha = 0.5) +
  facet_grid(duty_shift~shift_team) +
  scale_x_continuous(breaks = seq(2,12, by = 4)) +
  theme_minimal()
```



(You can remove this image once you've made yours.)

After discussing your investigation and the above graph with your Personnel Officer, they suggest you should *not* include rank, position, division, sub-division or gender in your analysis. They also tell you that ship-to-ship, how duty shifts are set up and how teams are allocated differs quite a lot. Some ships have more than the 4 shifts yours does, or have many more teams due to size, etc.

You're interested in presenting your work at the next Federation Science and Innovation Conference and want to be able to provide information that might be relevant to the Chief Science Officers on other ships, too.

Below are several models that you've fit and some tests on them.

```
model_1a <- lmer(productivity ~ week + starfleet_gpa + perseverance_score +
  (1|name),
  data = crew_data)

model_1b <- lmer(productivity ~ week + starfleet_gpa + perseverance_score +
  (1 + week|name),
```

```

data = crew_data)

# Study prompt: How do we interpret the p-values here? What is relevant?
lmtest::lrtest(model_1a, model_1b)

model_2a <- lmer(productivity ~ week + starfleet_gpa + perseverance_score +
  (1 + week|name) + (1|duty_shift:shift_team),
  data = crew_data)

model_2b <- lmer(productivity ~ week + starfleet_gpa + perseverance_score +
  (1 + week|name) + (1|full_team),
  data = crew_data)

# Study prompt: How do we interpret the p-values here? What is relevant?
lmtest::lrtest(model_1b, model_2a)
lmtest::lrtest(model_2a, model_2b)

```

2. Determine which model from the above is the most appropriate out of those shown. Make appropriate alterations to model_3 so that it will be the same as your chosen model with the addition of the term shown below, and uses the appropriate likelihood method to allow you to compare the models.

```

model_3 <- lmer(_____ +
  (1 + week|full_team),
  data = crew_data)

lmtest::lrtest(<your chosen model name here>, model_3)

```

3. Run `summary()` and `confint()` on whichever model you think is the most appropriate

Task set 5: competitive astrobiologists

While on shore leave, some of the astrobiologists had a little competition to see who could spot plants from the greatest number of **different planets or systems** in the hotel gardens. Note: The *number* of plants spotted doesn't actually matter as long as at least one was spotted.

They have asked for your impartial help to find out who the winner is.

You have three datasets:

- `astrobiologists` is a list of all the astrobiology crewmembers
- `competition_data` has the number of plants of each type that each participating astrobiologist recorded.
- `origin_data` contains information from the hotel about the plants in their collection and the the planets they are native to. They have warned you that is may be somewhat incomplete.

```
astrobiologists <- crew_data %>%
  filter(position == "Astrobiologist") %>%
  distinct(crew_id, name, .keep_all=TRUE) %>%
  transmute(crewmember = str_c(name, " (", crew_id, ")"))

competition_data <- tibble(crewmember =
  c(astrobiologists$crewmember[1],
    astrobiologists$crewmember[2],
    astrobiologists$crewmember[3]),
  `Xupta tree` = c(3L, 7L, NA),
  `L'maki` = c(21L, NA, 21L),
  `Leola root` = c(40L, 45L, 26L),
  Klavaatu = c(2L, 3L, 2L),
  Waterplum = c(NA, 5L, 1L),
  `Folnar jewel plant` = c(17L, 12L, 10L),
  `Felaran rose` = c(28L, 7L, NA),
  Crystilia = c(12L, 3L, 9L),
  Parthas = c(4L, 3L, NA),
  `Borgia plant` = c(NA, 1L, 1L))

origin_data <- data.frame(plant = c("Xupta tree", "L'maki", "Leola root",
  "Waterplum", "Vulcan orchid",
  "Lunar flower", "Garlanic tree",
  "Folnar jewel plant",
  "Felaran rose", "Crystilia", "Parthas",
  "Borgia plant", "Pod plant"),
  native_to = c("Orellius system", "Delta Quadrant",
  "Bajor", "Mari", "Vulcan",
  NA, "Elaysian homeworld", "Folnar III",
  "Delta Quadrant", "Telemarius IV",
  "Acomar III", "M-113", NA))
```

Tip: I recommend run `View()` on `competition_data` and `origin_data` to explore them further so you are familiar with their structure and contents. (You can also do this by clicking on their titles in the Environment pane.)

1. Create a new dataset called `complete_comp` using the `competition_data`.

2. Assess whether `complete_comp`, at this current step, is currently tidy. (I.e., is `competition_data` tidy?) If yes, proceed. If no, alter it to be tidy. Specifically, it needs to be in the correct format to be useful for merging the `origin_data` on to it.
3. Continuing to manipulate the `complete_comp` object, merge on the `origin_data` such that any plants **not** present in the data provided by the hotel are **dropped**.
4. Restrict the `complete_comp` so it only contains rows where at least one plant was spotted.
5. Restrict the `complete_comp` to just observations from distinct planets or systems for each crewmember. (See hint code below.)
6. Calculate how many unique planets or systems each astrobiologist spotted at least one plant from.

You DO NOT have to use the exact same code I do to get the associated questions in the timed component correct, as long as it fulfills these instructions, in the correct order. As a hint, here is the structure of my code to complete these tasks.

```
complete_comp <- competition_data %>%
  ----- %>%
  ----- %>%
  ----- %>%
distinct(crewmember, native_to) %>% # this line will achieve instruction 5
  ----- %>% # these last two lines achieve instruction 6,
  ----- # but could be done in only one line also
```

Task set 6: make your Statistician's life easier

Suppose you were trying to run the following code. It throws an error. (Note: DON'T fix the error, that isn't the point of this activity.) Create a reprex (a reproducible example, see week 1) with everything required for your statistician to reproduce this error. The only 'error' in the output should be the one produced by *this* code. (Hint: there is a library you should include, and you'll also need to provide the data. Once you've copied the complete code for the reprex to your clipboard, you can then run `reprex()` and the content for your reprex will then be added to your clipboard, (i.e., with Ctrl+V or Cmd+V you can paste it.))

```
origin_data %>%  
  filter(nativeto == "Delta Quadrant")
```

END