

# STAT 4710: Homework 1

Name

Due: September 14, 2023 at 12:00pm

## Contents

<b>Instructions</b>	<b>1</b>
<b>Case study: Major League Baseball</b>	<b>2</b>
<b>1 Wrangle (35 points for correctness; 5 points for presentation)</b>	<b>3</b>
1.1 Import (5 points)	3
1.2 Tidy (15 points)	3
1.3 Quality control (15 points)	3
<b>2 Explore (50 points for correctness; 10 points for presentation)</b>	<b>3</b>
2.1 Payroll across years (15 points)	3
2.2 Win percentage across years (15 points)	4
2.3 Win percentage versus payroll (15 points)	4
2.4 Team efficiency (5 points)	4

## Instructions

### Materials and collaboration

The policy on allowed materials and collaboration is as stated on the Syllabus:

“Students are permitted to work together on homework assignments, but must write up and submit solutions individually. In particular, students may not copy each others’ solutions. Students may consult all course materials, textbooks, the internet, or AI tools (e.g. ChatGPT or GitHub Copilot) to complete their homework. Students may not use solutions to problems that may be available online and/or from past iterations of the course. For each homework, students must disclose all classmates with whom they collaborated, which AI tools they used, and how they used them. Failure to do so will result in a 5-point penalty.”

In accordance with this policy,

*Please disclose all classmates with whom you collaborated:*

*Please disclose which AI tools you used, and how you used them:*

**Failure to answer the above questions will result in a 5-point penalty.**

### Writeup

Use this document as a starting point for your writeup, adding your solutions after “**Solution**”. Add your R code using code chunks and add your text answers using **bold text**. Consult the [preparing reports guide](#) for guidance on compilation, creation of figures and tables, and presentation quality. In particular, if the

instructions ask you to “print a table”, you should use `kable`. If the instructions ask you to “print a tibble”, you should not use `kable` and instead print the tibble directly.

## Programming

The `tidyverse` paradigm for data visualization, manipulation, and wrangling is required. No points will be awarded for code written in base R.

## Grading

The point value for each problem sub-part is indicated. Additionally, the presentation quality of the solution for each problem (as exemplified by the guidelines in Section 4 of the [preparing reports guide](#) will be evaluated on a per-problem basis (e.g. in this homework, there are three problems). There are 100 points possible on this homework, 85 of which are for correctness and 15 of which are for presentation.

## Submission

Compile your writeup to PDF and submit to Gradescope.

## Case study: Major League Baseball

What is the relationship between payroll and wins among Major League Baseball (MLB) teams? In this homework, we’ll find out by wrangling, exploring, and modeling the dataset in `MLPayData_Total.csv`, which contains the winning records and the payroll data of all 30 MLB teams from 1998 to 2014.

The dataset has the following variables:

- `payroll`: total team payroll (in billions of dollars) over the 17-year period
- `avgwin`: the aggregated win percentage over the 17-year period
- `Team.name.2014`: the name of the team
- `p1998, …, p2014`: payroll for each year (in millions of dollars)
- `X1998, …, X2014`: number of wins for each year
- `X1998.pct, …, X2014.pct`: win percentage for each year

We’ll need to use the following R packages:

```
library(tidyverse) # tidyverse
library(ggrepel)   # for scatter plot point labels
library(kableExtra) # for printing tables
library(cowplot)  # for side by side plots
```

# 1 Wrangle (35 points for correctness; 5 points for presentation)

## 1.1 Import (5 points)

- Import the data into a `tibble` called `mlb_raw` and print it.
- How many rows and columns does the data have?
- Do the numbers of rows and columns in the data match up with the data description given above?

**Solution.**

## 1.2 Tidy (15 points)

The raw data are in a messy format: Some of the column names are hard to interpret, we have data from different years in the same row, and both year-by-year and aggregate data are present.

- Tidy the data into two separate `tibbles`: one called `mlb_aggregate` containing the aggregate data and another called `mlb_yearly` containing the year-by-year data. `mlb_aggregate` should contain columns named `team`, `payroll_aggregate`, `pct_wins_aggregate` and `mlb_yearly` should contain columns named `team`, `year`, `payroll`, `pct_wins`, `num_wins`. Comment your code to explain each step.
- Print these two `tibbles`. How many rows do `mlb_aggregate` and `mlb_yearly` contain, and why?

**Solution.**

## 1.3 Quality control (15 points)

It's always a good idea to check whether a dataset is internally consistent. In this case, we are given both aggregated and yearly data, so we can check whether these match. To this end, carry out the following steps:

- Create a new `tibble` called `mlb_aggregate_computed` based on aggregating the data in `mlb_yearly`, containing columns named `team`, `payroll_aggregate_computed`, and `pct_wins_aggregate_computed`.
- Ideally, `mlb_aggregate_computed` would match `mlb_aggregate`. To check whether this is the case, join these two `tibbles` into `mlb_aggregate_joined` (which should have five columns: `team`, `payroll_aggregate`, `pct_wins_aggregate`, `payroll_aggregate_computed`, and `pct_wins_aggregate_computed`.)
- Create scatter plots of `payroll_aggregate_computed` versus `payroll_aggregate` and `pct_wins_aggregate_computed` versus `pct_wins_aggregate`, including a 45° line in each. Display these scatter plots side by side, and comment on the relationship between the computed and provided aggregate statistics.

**Solution.**

# 2 Explore (50 points for correctness; 10 points for presentation)

Now that the data are in tidy format, we can explore them by producing visualizations and summary statistics.

## 2.1 Payroll across years (15 points)

- Plot `payroll` as a function of `year` for each of the 30 teams, faceting the plot by `team` and adding a red dashed horizontal line for the mean payroll across years of each team.
- Using `dplyr`, identify the three teams with the greatest `payroll_aggregate_computed`, and print a table of these teams and their `payroll_aggregate_computed`.
- Using `dplyr`, identify the three teams with the greatest percentage increase in payroll from 1998 to 2014 (call it `pct_increase`), and print a table of these teams along with `pct_increase` as well as their payroll figures from 1998 and 2014.

- How are the metrics `payroll_aggregate_computed` and `pct_increase` reflected in the plot above, and how can we see that the two sets of teams identified above are the top three in terms of these metrics?

**Solution.**

## 2.2 Win percentage across years (15 points)

- Plot `pct_wins` as a function of `year` for each of the 30 teams, faceting the plot by `team` and adding a red dashed horizontal line for the average `pct_wins` across years of each team.
- Using `dplyr`, identify the three teams with the greatest `pct_wins_aggregate_computed` and print a table of these teams along with `pct_wins_aggregate_computed`.
- Using `dplyr`, identify the three teams with the most erratic `pct_wins` across years (as measured by the standard deviation, call it `pct_wins_sd`) and print a table of these teams along with `pct_wins_sd`.
- How are the metrics `pct_wins_aggregate_computed` and `pct_wins_sd` reflected in the plot above, and how can we see that the two sets of teams identified above are the top three in terms of these metrics?

**Solution.**

## 2.3 Win percentage versus payroll (15 points)

Let us investigate the relationship between win percentage and payroll.

- Using `mlb_aggregate_computed`, create a scatter plot of `pct_wins_aggregate_computed` versus `payroll_aggregate_computed`, labeling each point with the team name using `geom_text_repel` from the `ggrepel` package.
- Is the relationship between `payroll` and `pct_wins` positive or negative? Is this what you would expect, and why?

**Solution.**

## 2.4 Team efficiency (5 points)

Define a team's *efficiency* as the ratio of the aggregate win percentage to the aggregate payroll—more efficient teams are those that win more with less money.

- Using `dplyr`, identify the three teams with the greatest efficiency based on `mlb_aggregate_computed`, and print a table of these teams along with their efficiency, as well as their `pct_wins_aggregate_computed` and `payroll_aggregate_computed`.
- In what sense do these three teams appear efficient in the previous plot?

Side note: The movie “[Moneyball](#)” portrays “Oakland A’s general manager Billy Beane’s successful attempt to assemble a baseball team on a lean budget by employing computer-generated analysis to acquire new players.”

**Solution.**