



# STEAM PLAYTIME ANALYSIS

Data Science · Mr. Merrick · October 22, 2025

*You've been hired by Valve as a junior data scientist to explore trends in Steam game playtime data. Your mission: ask clear questions, explore patterns, and use R to uncover simple, data-driven insights.*

## Overview

Using the **Steam Playtime Dataset**, conduct open-ended exploratory analysis. Use R to summarize, visualize, and describe patterns in how players spend their time across games.

## Guiding Research Questions (Choose, Combine, or Extend)

Pick at least **three** prompts—or design your own.

### Q1. Dataset Overview.

How many rows (observations) and variables (columns) are there? What are the mean, median, spread, and range of playtimes?

### Q2. Overall Playtime Distribution.

Visualize and describe the overall distribution of **Hours**. Is it symmetric or skewed? Do a few games dominate total time?

### Q3. Specific Game Distribution (Skyrim).

Focus on *The Elder Scrolls V: Skyrim* (or a game of your choosing). Visualize the distribution of **Hours** for Skyrim and describe the pattern.

### Q4. Playtime Across a Franchise (Call of Duty).

How do *Call of Duty* games differ in playtime distributions? Use string filtering (e.g., `grep1`) to gather all *Call of Duty* titles. Compare the distributions of **Hours** across those games.

### Q5. Top Games by Median Playtime.

Which games have the highest *median* playtime? Show the top 10–15 and include a distribution plot so variation is visible.

### Q6. Most Popular Games (by Player Count).

Which games appear most often in the dataset (i.e., have the most rows/users)? Plot the distribution of player counts for all games.

### Q7. Popularity vs. Playtime.

Create a scatter plot of player counts vs. median playtime. Does there appear to be an association between the variables?

### Q8. Extending Further.

Since there are only two variables in the dataset we can only explore how 'name' explains variation in playtime. Can you think of some other measurable variables that might explain the variation in playtime?

## Dataset Description and Context

The dataset you'll use summarizes how long Steam users have played different games. Each row represents a player–game pair, meaning a single user's playtime for one game.

- **Name:** Title of the video game (e.g., *The Elder Scrolls V: Skyrim*)
- **Hours:** Number of hours the player has spent on that game

This dataset is adapted from the public [Steam Playtime Dataset](#) on Kaggle, simplified for beginner analysis. It's designed to let you explore patterns of popularity, engagement, and variation in playtime across titles.

## Solution Key (R Script)

```

1 # =====
2 # Steam Playtime Analysis - Solution Script
3 # Dataset columns: Name, Hours
4 # =====
5
6 # --- Libraries ---
7 library(readr)
8 library(dplyr)
9 library(ggplot2)
10 library(janitor)
11 library(scales)
12 library(stringr)
13
14 # --- 1) Load & clean ---
15 # Put the CSV in your working directory. Rename if needed.
16 st <- read_csv("steam.csv", show_col_types = FALSE)
17 st <- clean_names(st) # -> name, hours
18
19 # -----
20 # Q1. Dataset Overview
21 # -----
22 print(dim(st))
23 print(names(st))
24
25 summary_table <- st %>%
26   summarise(
27     n_rows = n(),
28     n_variables = ncol(st),
29     mean_hours = mean(hours),
30     median_hours = median(hours),
31     sd_hours = sd(hours),
32     min_hours = min(hours),
33     p25_hours = quantile(hours, 0.25),
34     p75_hours = quantile(hours, 0.75),
35     max_hours = max(hours)
36   )
37 print(summary_table)
38
39 # -----
40 # Q2. Overall Playtime Distribution
41 # -----
42 # Linear-scale histogram
43 ggplot(st, aes(x = hours)) +
44   geom_histogram(bins = 30, color = "white", fill = "skyblue") +
45   scale_x_continuous(labels = comma) +
46   labs(
47     title = "Distribution of Game Playtime (Hours)",
48     x = "Hours Played",
49     y = "Number of Records"
50   )
51
52 # Log-scale histogram (skew-friendly)
53 ggplot(st, aes(x = hours)) +
54   geom_histogram(bins = 30, color = "white", fill = "orange") +

```

```

55 scale_x_log10(labels = comma) +
56 labs(
57   title = "Distribution of Game Playtime (Log Scale)",
58   x = "Hours Played (log10)",
59   y = "Number of Records"
60 )
61
62 # -----
63 # Q3. Skyrim: Specific Game Distribution
64 # -----
65 st_skyrim <- st %>%
66   filter(grepl("skyrim", name, ignore.case = TRUE))
67
68 if (nrow(st_skyrim) > 0) {
69   # Summary
70   print(st_skyrim %>% summarise(
71     n_players = n(),
72     med_hours = median(hours),
73     mean_hours = mean(hours),
74     p90_hours = quantile(hours, 0.90)
75   ))
76
77   # Density + rug
78   ggplot(st_skyrim, aes(x = hours)) +
79     geom_histogram(bins = 30, color = "white", fill = "mediumseagreen") +
80     scale_x_continuous(labels = comma) +
81     labs(
82       title = "Playtime Distribution - The Elder Scrolls V: Skyrim",
83       x = "Hours Played",
84       y = "Number of Players"
85     )
86
87 } else {
88   message("Skyrim not found in dataset by the simple filter.")
89 }
90
91 # -----
92 # Q4. Franchise Comparison - Call of Duty
93 # -----
94 # Use grepl / base R style as requested
95 cod_rows <- st %>% filter(grepl("call of duty", name, ignore.case=TRUE))
96
97 if (nrow(cod_rows) > 0) {
98   # Order titles by median hours to keep plots tidy
99   cod_summary <- cod_rows %>%
100     group_by(name) %>%
101     summarise(
102       n_players = n(),
103       median_hours = median(hours),
104       mean_hours = mean(hours)
105     ) %>%
106     arrange(desc(median_hours))
107
108   print(cod_summary)
109
110   # Boxplots of hours per title (log scale to show spread)

```

```

111 ggplot(cod_rows,
112       aes(x = reorder(name, hours, FUN = median), y = hours)) +
113   geom_boxplot(fill = "lightblue", outlier.alpha = 0.3) +
114   coord_flip() +
115   labs(
116     title = "Call of Duty: Playtime Distributions by Title",
117     x = "Title",
118     y = "Hours Played"
119   )
120 } else {
121   message("No Call of Duty titles matched with grepl('call of duty', ...).")
122 }
123
124 # -----
125 # Q5. Top Games by Median Playtime (+ variation)
126 # -----
127 top_n <- 15
128 game_summary <- st %>%
129   group_by(name) %>%
130   summarise(
131     n_players = n(),
132     median_hours = median(hours),
133     mean_hours = mean(hours)
134   ) %>%
135   filter(n_players > 20) %>%
136   arrange(desc(median_hours))
137
138 top_games <- game_summary %>% slice_head(n = top_n)
139 print(top_games)
140
141 # Bar chart of medians
142 top_games %>% ggplot(aes(x = reorder(name, median_hours), y = median_hours)) +
143   geom_col(fill = "steelblue") +
144   coord_flip() +
145   scale_y_continuous(labels = comma) +
146   labs(
147     title = paste0("Top ", top_n, " Games by Median Playtime"),
148     x = "Game",
149     y = "Median Hours"
150   )
151
152 # Boxplots for variation among the same top games
153 st_top <- st %>% filter(name %in% top_games$name)
154 ggplot(st_top, aes(x = reorder(name, hours, FUN = median), y = hours)) +
155   geom_boxplot(fill = "lightgreen", outlier.alpha = 0.3) +
156   coord_flip() +
157   scale_y_log10(labels = comma) +
158   labs(
159     title = paste0("Playtime Distributions for Top ", top_n, " Games"),
160     x = "Game",
161     y = "Hours Played (log scale)"
162   )
163
164 # -----
165 # Q6. Most Popular Games (by player count)
166 # -----

```

```
167 popularity <- st %>%
168   group_by(name) %>%
169   summarize(player_count = n()) %>%
170   filter(player_count>20)
171
172 print(head(popularity, 15))
173
174 # Distribution of player counts across all games
175 ggplot(popularity, aes(x = player_count)) +
176   geom_histogram(bins = 30, color = "white", fill = "plum") +
177   scale_x_continuous(labels = comma) +
178   labs(
179     title = "Distribution of Player Counts Across Games",
180     x = "Player Count (rows per game)",
181     y = "Number of Games"
182   )
183
184 # -----
185 # Q7. Popularity vs Median Playtime (scatter)
186 # -----
187 pop_vs_play <- st %>%
188   group_by(name) %>%
189   summarise(
190     player_count = n(),
191     median_hours = median(hours)
192   )
193
194 ggplot(pop_vs_play, aes(x = player_count, y = median_hours)) +
195   geom_point(alpha = 0.3) +
196   labs(
197     title = "Popularity vs Median Playtime (per Game)",
198     x = "Player Count (log10)",
199     y = "Median Hours (log10)"
200   )
201
202 # End of script
```