



POKÉMON ANALYSIS

Data Science · Mr. Merrick · October 8, 2025

Dr. Oak has hired you as a junior data scientist to help him understand patterns and trends in Pokémon data. Your mission: ask questions, explore creatively, and use R to find insights worth sharing.

Overview

Using the Pokémon dataset, your task is to conduct open-ended exploratory research. You will use R to clean, visualize, and describe one-variable and simple two-variable patterns that help answer questions about Pokémon. A detailed description of the dataset can be found on page 2.

Guiding Research Questions (Choose, Combine, or Extend)

Dr. Oak encourages creativity. Use these as inspiration—but feel free to pose your own related questions. There are no “right answers,” only thoughtful, data-driven ones.

Q1. Summarize the Dataset.

How many observations? How many variables?

Q2. What does it mean for a Pokémon to be “strong” in Pokémon GO?

Define “strength” however you wish (e.g., Combat Power, Attack, or a combination). Which Pokémon best embody that strength? How common or rare are they?

Q3. How do different Pokémon *types* vary in their attributes?

Do certain types tend to be heavier, faster, or have higher CP? Are there clear patterns when you look across type distributions?

Q4. Are Legendary Pokémon noticeably stronger than non-Legendaries?

How do their stat distributions compare? Are they uniformly dominant, or do some normal Pokémon rival them?

Q5. Is there an association between a Pokémon’s weight and height?

Does size predict other characteristics like CP or speed? Are certain types (e.g., Flying, Water) exceptions to general trends?

Q6. Have newer generations of Pokémon become stronger or weaker?

Compare distributions of total strength or CP across generations. What trends emerge as new Pokémon are introduced?

Your report should explore **at least three** of these questions in depth and be submitted as a pdf.

Dataset Description and Context

The **Complete Pokémon Dataset** compiled by [Rounak Banik \(2018\)](#) is published on Kaggle. This dataset consolidates information from the official Pokémon main-series video games (*Pokémon Red/Blue* through *Pokémon Sun/Moon*, Generations I–VII). It includes one row per Pokémon species or form, with variables describing type, base statistics, generation, and legendary status.

- **Citation:** Banik, R. (2018). *The Complete Pokémon Dataset*. Kaggle. Retrieved from <https://www.kaggle.com/datasets/rounakbanik/pokemon>

These base statistics are drawn from Game Freak’s canonical data tables and are consistent with the in-game Pokédex entries from the core RPG series. No player-specific modifications (e.g., individual values or effort values) are included—these are species-level attributes.

Relation to Pokémon GO and the Trading Card Game

While all three media (the main games, Pokémon GO, and the Trading Card Game) share the same Pokémon species, they use different numerical systems for attributes:

System	Stats Used	Range / Scale	Notes
Main Games	HP, Attack, Defense, Special Attack, Special Defense, Speed	1–255	Turn-based RPG; 6 attributes per Pokémon.
Pokémon GO	Attack, Defense, Stamina	~10–300	Simplified for mobile play; 3 attributes + CP.
Trading Card Game (TCG)	HP, Attack Power (per move)	30–340 HP typical	Designed for card game balance.

Table 1: Comparison of stat systems across Pokémon media.

In the main series, base stats form the foundation of a deeper system involving:

- **Individual Values (IVs):** Random 0–31 bonuses unique to each Pokémon
- **Effort Values (EVs):** Training-based stat adjustments
- **Natures:** Personality modifiers that boost or reduce specific stats

Together, these create large individual variation between Pokémon of the same species—an important concept in probability and variation modeling.

By contrast, **Pokémon GO** collapses these six core stats into three (Attack, Defense, Stamina) to suit a real-time mobile environment. Individual variation is introduced via simplified IVs ranging from 0–15 per stat. The game’s “Combat Power (CP)” is then derived from these stats through a multiplicative formula.

The **Trading Card Game (TCG)** discards these mechanics entirely, using abstract HP and attack values chosen for card game balance rather than biological realism or direct translation from the video games.

Pokémon Data Science Tutorial — Group Tutorial (R)

Setup and Data Import

We'll load a few tidy packages, import the CSV, and clean column names.

```
1 # --- Libraries ---
2 library(readr)
3 library(dplyr)
4 library(ggplot2)
5 library(janitor)
6
7 # --- Load and clean dataset ---
8 pk <- read_csv("pokemon.csv") # ensure this file is in your working directory
9 pk <- clean_names(pk) # e.g., "Sp. Atk" -> "sp_atk"
10
11 # --- Preview the data ---
12 glimpse(pk)
13 summary(pk)
```

Q1. Summarize the Dataset

Prompt: How many observations (rows)? How many variables (columns)? What are the basic characteristics of the data?

```
1 # Rows and columns
2 dim(pk)
3
4 # Variable names
5 names(pk)
6
7 # --- Create a total stat variable (sum of base stats) ---
8 # NOTE: Depending on your CSV, the columns are typically: hp, attack, defense, sp_atk, sp_def,
9 #       speed.
10 # If your file uses "sp_attack"/"sp_defense", rename them or adjust the line below.
11 pk <- pk %>%
12   mutate(total = hp + attack + defense + sp_atk + sp_def + speed)
13
14 # Simple summary statistics
15 pk %>%
16   summarise(
17     n_pokemon = n(),
18     n_variables = ncol(pk),
19     mean_total = mean(total, na.rm = TRUE),
20     median_total = median(total, na.rm = TRUE)
21   )
22
23 # Average total stats by Generation
24 pk %>%
25   group_by(generation) %>%
26   summarise(avg_total = mean(total, na.rm = TRUE))
```

Q2. What Does it Mean for a Pokémon to be “Strong”?

Prompt: We'll define strength as the sum of the six base stats (`total`). Find the strongest Pokémon overall and visualize the distribution.

```
1 # Top 10 Pokemon by total stats
2 pk %>%
3   arrange(desc(total)) %>%
4   select(name, type_1, type_2, total) %>%
5   head(10)
6
7 # Simple histogram of total strength
8 ggplot(pk, aes(x = total)) +
9   geom_histogram(bins = 25, fill = "skyblue", color = "white") +
10  labs(title = "Distribution of Total Pokemon Strength",
11        x = "Total Base Stats", y = "Count")
```

Q3. How Do Different Types Vary?

Prompt: Compare average strength across primary types and show a type comparison plot.

```
1 # Average total stats by primary type
2 pk %>%
3   group_by(type_1) %>%
4   summarise(avg_total = mean(total, na.rm = TRUE)) %>%
5   arrange(desc(avg_total))
6
7 # Boxplot comparing types
8 ggplot(pk, aes(x = type_1, y = total)) +
9   geom_boxplot(fill = "lightgreen") +
10  coord_flip() +
11  labs(title = "Total Stats by Primary Type", x = "Type", y = "Total Stats")
```

Q4. Are Legendary Pokémon Stronger?

Prompt: Compare totals for Legendary vs. non-Legendary. Convert the 0/1 indicator to labeled categories for plotting.

```
1 # Summary table by legendary indicator (0/1)
2 pk %>%
3   group_by(is_legendary) %>%
4   summarise(avg_total = mean(total, na.rm = TRUE),
5             count = n())
6
7 # Make a labeled factor for plotting
8 pk <- pk %>%
9   mutate(is_legendary_factored = factor(is_legendary,
10                                         levels = c(0, 1),
11                                         labels = c("No", "Yes")))
12
13 # Boxplot: Legendary vs Non-Legendary
14 ggplot(pk, aes(x = is_legendary_factored, y = total, fill = is_legendary_factored)) +
15   geom_boxplot() +
16   labs(title = "Legendary vs Non-Legendary Pokemon",
17        x = "Legendary?",
18        y = "Total Stats") +
```

```

19 scale_fill_manual(values = c("No" = "skyblue", "Yes" = "orange"))
20
21 # Faceted histograms (same y-scale)
22 ggplot(pk, aes(x = total)) +
23   geom_histogram(bins = 25, color = "white", fill = "skyblue") +
24   facet_wrap(~ isLegendaryFactored, ncol = 2) +
25   labs(title = "Distribution of Total Base Stats",
26        x = "Total (BST)",
27        y = "Count")
28
29 # Faceted histograms (free y-axis to handle different group sizes)
30 ggplot(pk, aes(x = total)) +
31   geom_histogram(bins = 25, color = "white", fill = "skyblue") +
32   facet_wrap(~ isLegendaryFactored, ncol = 2, scales = "free_y") +
33   labs(title = "Distribution of Total Base Stats",
34        x = "Total (BST)",
35        y = "Count")

```

Q5. Is There an Association Between Height and Weight?

Prompt: Make a simple scatter plot and (optionally) compute correlation.

```

1 ggplot(pk, aes(x = height_m, y = weight_kg)) +
2   geom_point(alpha = 0.7, color = "darkblue") +
3   labs(title = "Height vs Weight of Pokemon", x = "Height (m)", y = "Weight (kg)")
4
5 # Optional: correlation
6 cor(pk$height_m, pk$weight_kg, use = "complete.obs")

```

Q6. Have Newer Generations Become Stronger or Weaker?

Prompt: Compare total strength across generations.

```

1 ggplot(pk, aes(x = factor(generation), y = total, fill = factor(generation))) +
2   geom_boxplot() +
3   labs(title = "Pokemon Strength Across Generations",
4        x = "Generation", y = "Total Stats") +
5   scale_fill_brewer(palette = "Set3")
6
7 pk %>%
8   group_by(generation) %>%
9   summarise(avg_total = mean(total, na.rm = TRUE),
10            median_total = median(total, na.rm = TRUE))

```

Wrap-Up

We used simple summaries and visuals to explore patterns: overall strength, type differences, legendary comparisons, size relationships, and generation trends. In future assignments, you'll extend these ideas and write short interpretations for each figure or table.

Python Supplement (pandas + seaborn)

Setup & Data Import

```

1 # If needed:
2 # !pip install pandas seaborn matplotlib
3
4 import pandas as pd
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7
8 sns.set_theme(context="notebook", style="whitegrid")
9
10 # Load CSV (same directory as your notebook / script)
11 pk = pd.read_csv("pokemon.csv")
12
13 print(pk.shape)
14 print(pk.columns.tolist())

```

Q1. Summarize the Dataset

```

1 # Create 'total' (Base Stat Total) if you have these columns
2 stat_cols = ["hp", "attack", "defense", "sp_atk", "sp_def", "speed"]
3 missing = [c for c in stat_cols if c not in pk.columns]
4 if missing:
5     print("Missing:", missing)
6 else:
7     pk["total"] = pk[stat_cols].sum(axis=1)
8
9 # Rows, columns
10 n_rows, n_cols = pk.shape
11 print("observations:", n_rows, "| variables:", n_cols)
12
13 # Simple summary
14 if "total" in pk.columns:
15     print(pk["total"].describe()[["mean", "50%"]]) # mean, median
16
17 # Average total by Generation
18 if {"generation", "total"}.issubset(pk.columns):
19     print(pk.groupby("generation")["total"].mean())

```

Q2. What Does it Mean to be “Strong”?

```

1 # Define strength as 'total' and list strongest Pok\emon
2 cols = [c for c in ["name", "type_1", "type_2", "total"] if c in pk.columns]
3 if "total" in pk.columns and cols:
4     print(pk.sort_values("total", ascending=False)[cols].head(10))
5
6 # Histogram of total
7 if "total" in pk.columns:
8     plt.figure()
9     sns.histplot(pk["total"], bins=25)
10    plt.title("Distribution of Total Pok\emon Strength")

```

```

11 plt.xlabel("Total Base Stats"); plt.ylabel("Count")
12 plt.tight_layout(); plt.show()

```

Q3. How Do Types Vary?

```

1 # Average total by primary type
2 if {"type_1", "total"}.issubset(pk.columns):
3     avg_by_type = (pk.groupby("type_1")["total"]
4                     .mean().reset_index(name="avg_total")
5                     .sort_values("avg_total", ascending=False))
6     print(avg_by_type)
7
8 # Boxplot by type
9 if {"type_1", "total"}.issubset(pk.columns):
10    plt.figure(figsize=(8,6))
11    sns.boxplot(data=pk, x="type_1", y="total")
12    plt.title("Total Stats by Primary Type")
13    plt.xlabel("Type"); plt.ylabel("Total Stats")
14    plt.xticks(rotation=45, ha="right")
15    plt.tight_layout(); plt.show()

```

Q4. Are Legendary Pokémon Stronger?

```

1 # Summaries by legendary status
2 if {"is_legendary", "total"}.issubset(pk.columns):
3     print(pk.groupby("is_legendary")["total"]
4           .agg(avg_total="mean", count="size")
5           .reset_index())
6
7 # Make a labeled flag for plotting
8 if pk["is_legendary"].dtype != "0":
9     pk["is_legendary_flag"] = pk["is_legendary"].astype(int)
10 else:
11     pk["is_legendary_flag"] = pk["is_legendary"].astype(str).str.lower().map({
12         "0": "0", "1": "1", "false": "0", "true": "1", "no": "0", "yes": "1"
13     }).fillna("0").astype(int)
14 pk["is_legendary_label"] = pk["is_legendary_flag"].map({0: "No", 1: "Yes"})
15
16 # Boxplot
17 plt.figure()
18 sns.boxplot(data=pk, x="is_legendary_label", y="total")
19 plt.title("Legendary vs Non-Legendary (Total Stats)")
20 plt.xlabel("Legendary?"); plt.ylabel("Total Stats")
21 plt.tight_layout(); plt.show()
22
23 # Faceted histograms (same y)
24 g = sns.FacetGrid(pk, col="is_legendary_label", sharey=True)
25 g.map_dataframe(sns.histplot, x="total", bins=25)
26 g.set_axis_labels("Total (BST)", "Count")
27 g.fig.subplots_adjust(top=0.85)
28 g.fig.suptitle("Distribution of Total Base Stats")
29 plt.show()
30
31 # Faceted histograms (free y)

```

```
32 g2 = sns.FacetGrid(pk, col="is_legendary_label", sharey=False)
33 g2.map_dataframe(sns.histplot, x="total", bins=25)
34 g2.set_axis_labels("Total (BST)", "Count")
35 g2.fig.subplots_adjust(top=0.85)
36 g2.fig.suptitle("Distribution of Total Base Stats (free y)")
37 plt.show()
```

Q5. Height vs Weight

```
1 if {"height_m", "weight_kg"}.issubset(pk.columns):
2     plt.figure()
3     sns.scatterplot(data=pk, x="height_m", y="weight_kg", alpha=0.7)
4     plt.title("Height vs Weight of Pok\`emon")
5     plt.xlabel("Height (m)"); plt.ylabel("Weight (kg)")
6     plt.tight_layout(); plt.show()
7
8     print("Correlation:",
9           pk[["height_m", "weight_kg"]].corr().iloc[0,1])
```

Q6. Generations: Stronger or Weaker?

```
1 if {"generation", "total"}.issubset(pk.columns):
2     plt.figure()
3     sns.boxplot(data=pk, x=pk["generation"].astype("category"), y="total")
4     plt.title("Pok\`emon Strength Across Generations")
5     plt.xlabel("Generation"); plt.ylabel("Total Stats")
6     plt.tight_layout(); plt.show()
7
8     print(pk.groupby("generation")["total"]
9           .agg(avg_total="mean", median_total="median")
10            .reset_index())
```