



Explore how LEGO set size, price, themes, and release year interact. Your goal: simple summaries, clean plots, and short interpretations.

Overview

Using the **LEGO Sets and Prices Over Time** dataset, conduct a light exploratory analysis in **R**. Include 2–3 sentence descriptions under each figure/table. Cover at least **three** questions (Q2 on missing data counts as one).

Guiding Questions

Q1. Summarize the dataset.

Rows, columns, year range, average pieces and average price. Describe the distribution of set piece counts.

Q2. Missing data overview.

Which variables (`name`, `theme`, `year`, `pieces`, `usd_msrp`) have missing values? Show a percent-missing bar chart and state how you will handle missingness in later questions.

Q3. What is “value”? (Price per piece)

Define $PPP = \text{usd_msrp}/\text{pieces}$. List several “best value” sets (e.g., ≥ 100 pieces). Show a PPP histogram. *Note: PPP can be very skewed; you may filter clear outliers (e.g., $PPP < \$1$). Say what you did and comment on why the distribution of PPP is heavily skewed.*

Q4. How do themes differ?

Compare themes by average PPP and/or average pieces. Limit to themes with at least 10 sets. Show a table or bar chart.

Q5. Two-theme comparison (pick any two).

Compare two specific themes (e.g., *Star Wars* vs *Harry Potter*) on price, pieces, and PPP. Include a small summary table and one plot (e.g., two boxplots or a side-by-side bar chart). Briefly interpret the differences.

Q6. Are bigger sets always pricier?

Make a scatter plot of `pieces` vs `usd_msrp` with a trend line and report the correlation.

Q7. How have sets changed over time?

Plot yearly averages for pieces, price, and/or PPP and describe any trends. Also explore the average changes in price/pieces/PPP for Star Wars, Disney, and Harry Potter themed sets.

Dataset Description and Citation

The dataset includes set name/ID, theme, year, piece counts, and prices.

- **Citation:** Alex Racapé. *LEGO Sets and Prices Over Time*. Kaggle. Retrieved from <https://www.kaggle.com/datasets/alexracape/lego-sets-and-prices-over-time/data>

Group Tutorial (R only)

Setup and Data Import

```

1 # =====
2 # Intro to R: LEGO Sets (Prices & Pieces)
3 # Gentle EDA + Missing Data basics (R only)
4 # Packages: readr, dplyr, ggplot2, janitor
5 # File needs: name, theme, year, pieces, usd_msrp
6 # =====
7 library(readr)
8 library(dplyr)
9 library(ggplot2)
10 library(janitor)
11
12 # Read & clean -----
13 lego <- read_csv("lego.csv", show_col_types = FALSE) |>
14   clean_names()
15
16 # If needed (uncomment):
17 # lego <- lego |>
18 # rename(pieces = num_parts, usd_msrp = retail_price)
19
20 names(lego)
21 dim(lego)
22 head(lego, 5)

```

Q1. Summarize

```

1 lego |>
2   summarise(
3     rows = n(),
4     cols = ncol(lego),
5     avg_pieces = mean(pieces, na.rm = TRUE),
6     avg_price = mean(usd_msrp, na.rm = TRUE),
7     first_year = min(year, na.rm = TRUE),
8     last_year = max(year, na.rm = TRUE)
9   )
10
11 ggplot(lego, aes(x = pieces)) +
12   geom_histogram(bins = 30, color = "white", fill = "steelblue") +
13   labs(title = "Distribution of LEGO Set Piece Counts",
14        x = "Pieces", y = "Count")

```

Q2. Missing Data Overview

```

1 pct <- function(x) round(100 * x, 1)
2
3 missing_overview <- lego |>
4   summarise(
5     n = n(),
6     missing_name = sum(is.na(name)),
7     missing_theme = sum(is.na(theme)),

```

```

8   missing_year = sum(is.na(year)),
9   missing_pieces = sum(is.na(pieces)),
10  missing_usd = sum(is.na(usd_msrp))
11 ) |>
12 mutate(
13   pct_name = pct(missing_name / n),
14   pct_theme = pct(missing_theme / n),
15   pct_year = pct(missing_year / n),
16   pct_pieces = pct(missing_pieces / n),
17   pct_usd = pct(missing_usd / n)
18 )
19
20 print(missing_overview)
21
22 miss_long <- tibble(
23   variable = c("name", "theme", "year", "pieces", "usd_msrp"),
24   pct_missing = c(
25     pct(mean(is.na(lego$name))),
26     pct(mean(is.na(lego$theme))),
27     pct(mean(is.na(lego$year))),
28     pct(mean(is.na(lego$pieces))),
29     pct(mean(is.na(lego$usd_msrp)))
30   )
31 )
32
33 ggplot(miss_long, aes(x = variable, y = pct_missing)) +
34   geom_col(fill = "tomato") +
35   labs(title = "Percent Missing by Variable",
36        x = "Variable", y = "Percent missing") +
37   ylim(0, 100)

```

Q3. Value (PPP)

```

1  # Naive definition (for discussion)
2  # lego |> mutate(ppp = usd_msrp/pieces)
3
4  # Safe PPP with guards for zeros / NAs
5  lego <- lego |>
6    mutate(ppp = ifelse(!is.na(usd_msrp) & !is.na(pieces) & pieces > 0,
7                      usd_msrp / pieces, NA_real_))
8
9  lego |>
10   filter(pieces >= 100, is.finite(ppp)) |>
11   arrange(ppp) |>
12   select(name, theme, year, pieces, usd_msrp, ppp) |>
13   print(n = 50)
14
15 # Distribution including outliers
16 ggplot(lego |> filter(!is.na(ppp)), aes(x = ppp)) +
17   geom_histogram(bins = 30, color = "white", fill = "seagreen") +
18   labs(title = "Distribution of Price per Piece (PPP)",
19        x = "USD per piece", y = "Count")
20
21 # Inspect high PPP entries
22 lego |>

```

```

23 arrange(desc(ppp)) |>
24 select(name, theme, year, pieces, usd_msrp, ppp) |>
25 print(n = 100)
26
27 # Optional: trimmed view for a cleaner shape
28 lego_trim <- lego |> filter(!is.na(ppp), ppp > 0, ppp < 1.5)
29 ggplot(lego_trim, aes(x = ppp)) +
30   geom_histogram(bins = 30, color = "white", fill = "seagreen") +
31   labs(title = "PPP (filtered < $1.5)",
32         x = "USD per piece", y = "Count")

```

Q4. Themes Compared

```

1 theme_summary <- lego |>
2   filter(!is.na(ppp)) |>
3   group_by(theme) |>
4   summarise(
5     sets = n(),
6     avg_ppp = mean(ppp, na.rm = TRUE),
7     avg_pieces = mean(pieces, na.rm = TRUE),
8     .groups = "drop"
9   ) |>
10  filter(sets >= 10) |>
11  arrange(avg_ppp)
12
13 print(theme_summary, n = 100)
14
15 top40 <- theme_summary |> slice_min(avg_ppp, n = 40)
16
17 ggplot(top40, aes(x = reorder(theme, avg_ppp), y = avg_ppp)) +
18   geom_col(fill = "tan3") +
19   coord_flip() +
20   labs(title = "Top Themes by Value (Lowest Avg PPP)",
21         x = "Theme", y = "Average USD per piece")
22
23 # (All-themes boxplot often messy due to many categories)
24 lego |> filter(!is.na(ppp)) |>
25   ggplot(aes(x = theme, y = usd_msrp)) + geom_boxplot()

```

Q5. Two-Theme Comparison

```

1 themes_to_compare <- c("Star Wars", "Harry Potter") # edit as desired
2
3 pair <- lego |> filter(theme %in% themes_to_compare)
4
5 pair |> group_by(theme) |>
6   summarise(
7     sets = n(),
8     avg_price = mean(usd_msrp, na.rm = TRUE),
9     avg_pieces = mean(pieces, na.rm = TRUE),
10    avg_ppp = mean(ppp, na.rm = TRUE),
11    .groups = "drop"
12  )
13

```

```

14 ggplot(pair |> filter(!is.na(usd_msrp)),
15       aes(x = theme, y = usd_msrp)) +
16   geom_boxplot(fill = "skyblue") +
17   labs(title = "Price by Theme (choose two)",
18       x = "Theme", y = "Price (USD)")

```

Q6. Pieces vs Price

```

1 lego_complete_price <- lego |>
2   filter(!is.na(pieces), !is.na(usd_msrp))
3
4 ggplot(lego_complete_price, aes(x = pieces, y = usd_msrp)) +
5   geom_point(alpha = 0.5) +
6   geom_smooth(method = "lm", se = TRUE) +
7   labs(title = "Pieces vs. Price",
8       x = "Pieces", y = "Price (USD)")
9
10 cor(lego_complete_price$pieces, lego_complete_price$usd_msrp)

```

Q7. Trends Over Time

```

1 yearly <- lego |>
2   group_by(year) |>
3   summarise(
4     avg_pieces = mean(pieces, na.rm = TRUE),
5     avg_price = mean(usd_msrp, na.rm = TRUE),
6     avg_ppp = mean(ppp, na.rm = TRUE),
7     sets = n()
8   ) |>
9   filter(!is.na(year))
10
11 ggplot(yearly, aes(x = year, y = avg_pieces)) +
12   geom_line() + geom_point() +
13   labs(title = "Average Pieces by Year",
14       x = "Year", y = "Average pieces")
15
16 ggplot(yearly, aes(x = year, y = avg_price)) +
17   geom_line() + geom_point() +
18   labs(title = "Average Price by Year",
19       x = "Year", y = "Average USD")
20
21 ggplot(yearly, aes(x = year, y = avg_ppp)) +
22   geom_line() + geom_point() +
23   labs(title = "Average Price per Piece (PPP) by Year",
24       x = "Year", y = "Average USD per piece")

```

Themed Lines Over Time

```

1 yearly <- lego |>
2   filter(theme %in% c("Star Wars", "Harry Potter", "Disney")) |>
3   group_by(year, theme) |>
4   summarise(

```

```
5     avg_pieces = mean(pieces, na.rm = TRUE),
6     avg_price = mean(usd_msrp, na.rm = TRUE),
7     avg_ppp = mean(ppp, na.rm = TRUE),
8     sets = n()
9   ) |>
10  filter(!is.na(year))
11
12  ggplot(yearly, aes(x = year, y = avg_pieces, colour = theme)) +
13    geom_line() + geom_point() +
14    labs(title = "Average Pieces by Year",
15         x = "Year", y = "Average pieces")
16
17  ggplot(yearly, aes(x = year, y = avg_price, colour = theme)) +
18    geom_line() + geom_point() +
19    labs(title = "Average Price by Year",
20         x = "Year", y = "Average USD")
21
22  ggplot(yearly, aes(x = year, y = avg_ppp, colour = theme)) +
23    geom_line() + geom_point() +
24    labs(title = "Average Price per Piece (PPP) by Year",
25         x = "Year", y = "Average USD per piece")
```