# Problem Set 3

## Data Visualisation for Social Scientists

### Due: February 18, 2026

## Instructions

- Please show your work! You may lose points by simply writing in the answer. If the problem requires you to execute commands in R, please include the code you used to get your answers. Please also include the .R file that contains your code. If you are not sure if work needs to be shown for a particular problem, please ask.

- Your homework should be submitted electronically on GitHub.

- This problem set is due before 23:59 on Wednesday February 18, 2026. No late assignments will be accepted.

## Canadian Election Study

The data for this problem set come from the Canadian Election Study (CES) in 2015. The main purpose of the study is to give a comprehensive picture of the Canadian election: why people vote as they do, what changes during campaigns and across elections, and how Canadian voting compares with that in other democracies.

### Data Manipulation

1. Load the CES .csv file from GitHub into your global environment. Filter respondents to only include "high quality" participants:

   ```
   ces2015 <- ces2015 |> filter(discard == "Good quality")
   ```

   ```
   1 # Load data
   2 ces2015 <- read.csv("CES2015.csv")
   3
   4 # Q1: Filter for "Good quality" participants
   5 ces2015 <- ces2015 %>% filter(discard == "Good quality")
   ```

2. Filter the dataset to those participants that answered the question about voting for the past election using `p_voted`. Consider respondents who gave a "Yes" answer as having voted, while "No" as not having voted. Treat "Don't know" and "Refused" as missing.

```
1  ces2015 <- ces2015 %>%
2    mutate(voted_binary = case_when(
3      p_voted == "Yes" ~ 1,
4      p_voted == "No" ~ 0,
5      TRUE ~ NA_real_
6    )) %>%
7    filter(!is.na(voted_binary))
```

3. Create an age variable and group into categories (e.g., <30, 30-44, 45-64, 65+). Year of birth is in age (four-digit year).

```
1  ces2015 <- ces2015 %>%
2    mutate(age = as.numeric(age)) %>%
3    mutate(actual_age = 2015 - age) %>%
4    mutate(age_group = case_when(
5      actual_age < 30 ~ "<30",
6      actual_age >= 30 & actual_age <= 44 ~ "30-44",
7      actual_age >= 45 & actual_age <= 64 ~ "45-64",
8      actual_age >= 65 ~ "65+",
9      TRUE ~ NA_character_
10   ))
```

## Data Visualization

1. Plot turnout rate by age group.

   **Answer:** I visualized the turnout rate by age group using a bar chart. A key step was converting the age categories into factors to ensure a logical chronological flow. This prevents R's default alphabetical sorting, which would mislead the reader. The Y-axis was formatted using `scales::percent` to clearly communicate the probability of voting in each cohort.
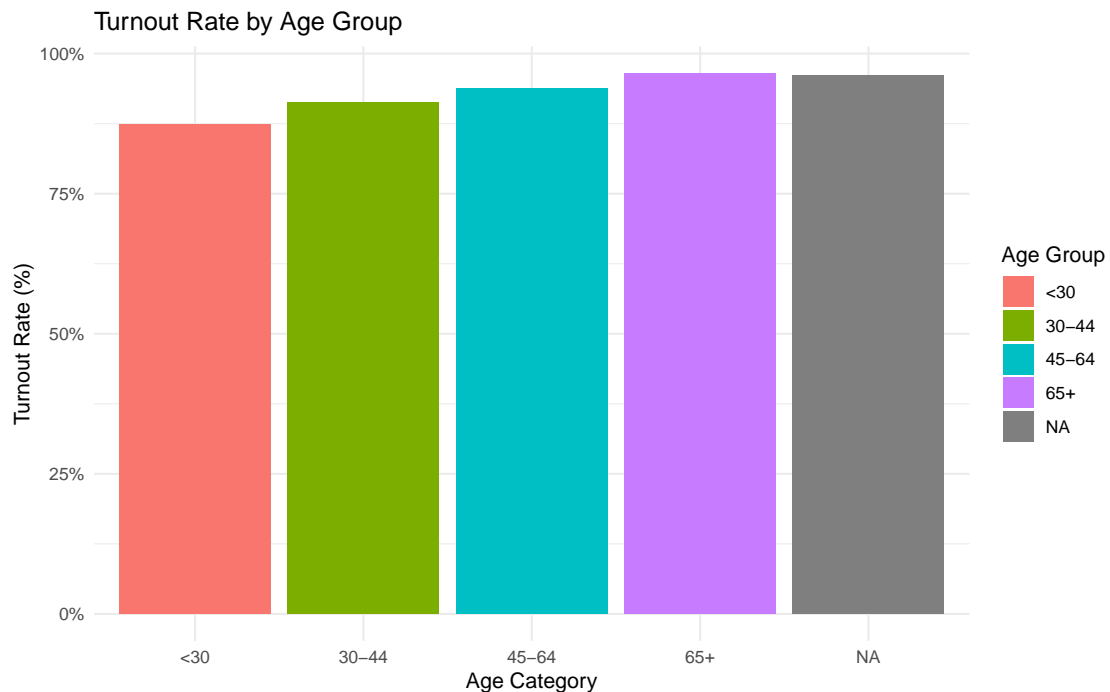
```
1  # 1. Define the logical chronological order for age groups
2  # This prevents the default alphabetical sorting (which would put <30 at
      the end).
3  age_levels <- c("<30", "30-44", "45-64", "65+")
4
5  plot1 <- ces2015 %>%
6    # 2. Convert age_group to a factor with specific levels
7    mutate(age_group = factor(age_group, levels = age_levels)) %>%
8
9    # 3. Aggregate data to calculate the mean turnout rate per group
10   group_by(age_group) %>%
11   summarise(turnout_rate = mean(voted_binary, na.rm = TRUE)) %>%
```

```
12
13    # 4. Initialize ggplot
14    ggplot(aes(x = age_group, y = turnout_rate, fill = age_group)) +
15
16    # Use geom_col for bar chart representation
17    geom_col() +
18
19    # Use scales::percent for the Y-axis to avoid 'object not found' errors
20    scale_y_continuous(labels = scales::percent) +
21
22    # 5. Add clear labels to enhance the plot's truthfulness [cite: 1847]
23    labs(
24        title = "Turnout Rate by Age Group",
25        x = "Age Category",
26        y = "Turnout Rate (%)",
27        fill = "Age Group"
28    ) +
29    theme_minimal()
```


Turnout Rate by Age Group

2. Create a density plot of ideology by party, restricting your sample to respondents with non-missing left–right self-placement (0–10 scale) and those that intended to vote for a main party (e.g., Liberal, Conservative, NDP, Bloc in Quebec, and Green).
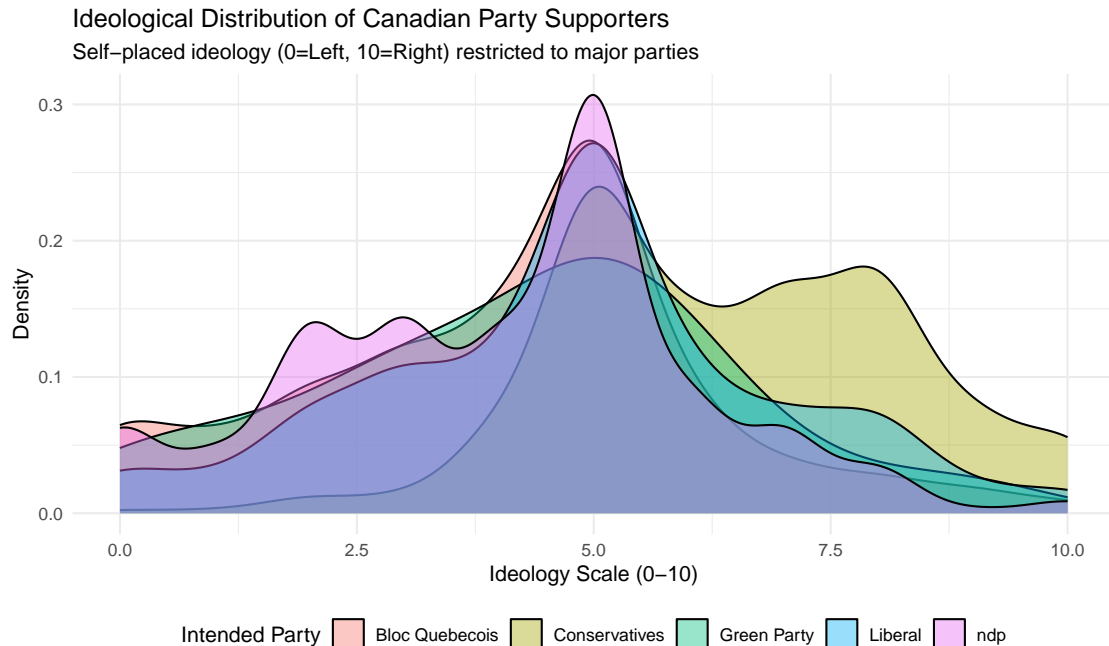
   **Answer:** This density plot shows the ideological distribution across major parties. I first performed data cleaning by removing non-substantive outliers like '1000' from the `p_selfplace` variable. I applied alpha transparency (0.4) to manage the overlapping

distributions, making it easier to see where supporters of different parties cluster on the 0-10 left-right scale.

```r
# 1. Update party list to match actual data values (Spelling & Case
    Sensitive)
# Data inspection revealed specific strings: "Conservatives" (plural), "
    ndp" (lowercase),
# and "Green Party" (full name).
main_parties <- c("Liberal", "Conservatives", "ndp", "Bloc Quebecois", "
    Green Party")

plot2 <- ces2015 %>%
# 2. Convert ideology variable (p_selfplace) to numeric and clean
# This fixes the "non-numeric argument" error by forcing text strings to
    NA.
mutate(p_selfplace = as.numeric(p_selfplace)) %>%

# 3. Filter for valid 0-10 range and major party supporters
# This step removes outliers (e.g., '1000') to focus on the substantive
    distribution.
filter(p_selfplace >= 0 & p_selfplace <= 10) %>%
filter(p_votechce %in% main_parties) %>%

# 4. Initialize ggplot visualization
ggplot(aes(x = p_selfplace, fill = p_votechce)) +

# Set transparency (alpha) to visualize overlapping distributions
geom_density(alpha = 0.4) +

# 5. Add labels and configure theme for professional reporting
labs(
  title = "Ideological Distribution of Canadian Party Supporters",
  subtitle = "Self-placed ideology (0=Left, 10=Right) restricted to major
      parties",
  x = "Ideology Scale (0-10)",
  y = "Density",
  fill = "Intended Party",
  caption = "Source: 2015 CES. Note: Only values 0-10 included; 'ndp' and
      'Green Party' labels matched to raw data."
) +

# Apply theme elements for better alignment and reduced non-data ink
theme_minimal() +
theme(legend.position = "bottom")

# Export the plot to PDF for the LaTeX report
ggsave("Plot2.pdf", plot = plot2, width = 8, height = 5)
```

**Ideological Distribution of Canadian Party Supporters**
Self–placed ideology (0=Left, 10=Right) restricted to major parties



Source: 2015 CES. Note: Only values 0–10 included; 'ndp' and 'Green Party' labels matched to raw data.

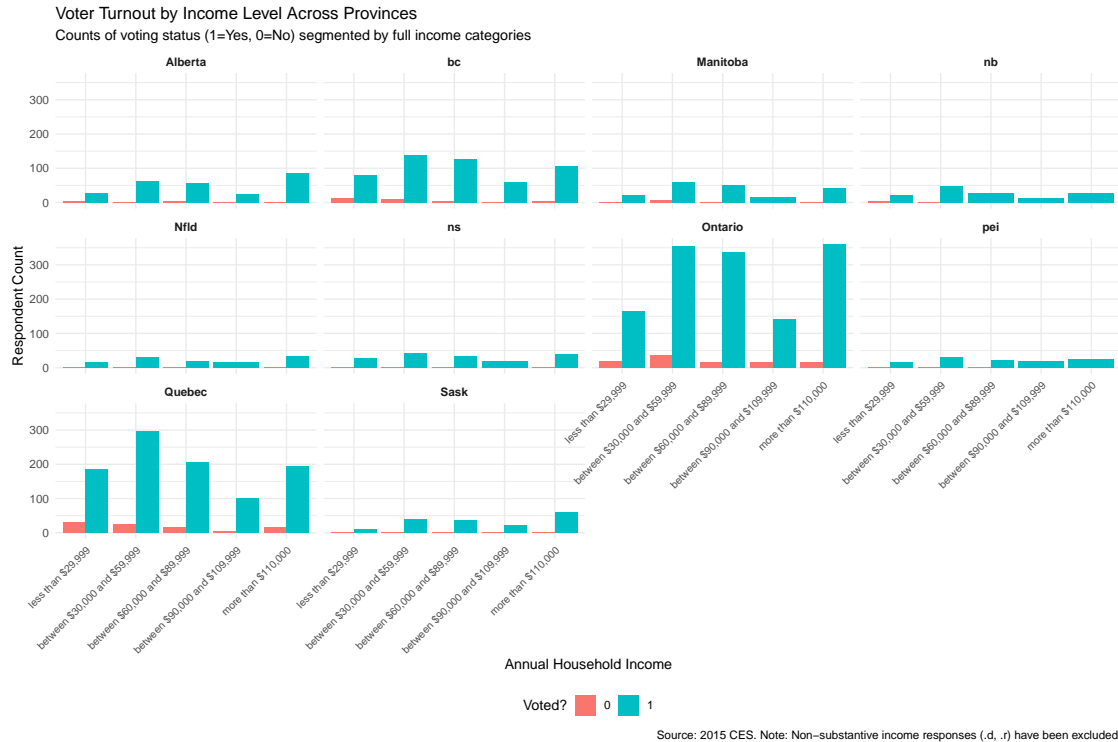3. Produce histogram counts of turnout by income (`income_full`), faceted by province.

   **Answer:** To analyze the relationship between income levels and turnout across different regions, I produced a histogram faceted by province. I utilized the **small multiples** technique via `facet_wrap()` to allow for easy geographic comparisons without overcrowding a single plot. The `income_full` variable was reordered into a logical progression from "less than $29,999" to "more than $110,000" to ensure the plot's **truthfulness** and intuitive readability.

```r
income_levels <- c("less than $29,999",
                   "between $30,000 and $59,999",
                   "between $60,000 and $89,999",
                   "between $90,000 and $109,999",
                   "more than $110,000")

plot3 <- ces2015 %>%
  # 2. Filter out non-substantive responses like '.r' or '.d'
  # and exclude invalid province codes like '1000'.
  filter(income_full %in% income_levels) %>%
  filter(province != "1000") %>%

  # 3. Convert income_full to a factor with the specified order
  # This follows the principle of "Consistent Ordering" from Week 4.
  mutate(income_full = factor(income_full, levels = income_levels)) %>%

  # 4. Initialize ggplot
```

```
18   # Using as.factor(voted_binary) to create discrete colors for Yes/No
19   ggplot(aes(x = income_full, fill = as.factor(voted_binary))) +
20
21   # Use geom_bar with position="dodge" for side-by-side count comparison
22   geom_bar(position = "dodge") +
23
24   # Implement Small Multiples by faceting by province
25   # This allows for easy geographic comparison as per instructions.
26   facet_wrap(~ province) +
27
28   # 5. Labels and appearance
29   theme_minimal() +
30   theme(
31     axis.text.x = element_text(angle = 45, hjust = 1, size = 8),
32     legend.position = "bottom",
33     strip.text = element_text(face = "bold") # Make province names bold
34   ) +
35   labs(
36     title = "Voter Turnout by Income Level Across Provinces",
37     subtitle = "Counts of voting status (1=Yes, 0=No) segmented by full
      income categories",
38     x = "Annual Household Income",
39     y = "Respondent Count",
40     fill = "Voted?",
41     caption = "Source: 2015 CES. Note: Non-substantive income responses
      (.d, .r) have been excluded."
42   )
```

Voter Turnout by Income Level Across Provinces
Counts of voting status (1=Yes, 0=No) segmented by full income categories

Respondent Count

Annual Household Income

Voted? 0 1

Source: 2015 CES. Note: Non–substantive income responses (.d, .r) have been excluded.

4. Create your own reusable custom theme. Apply your theme to one of the previous plots and add:

    (a) An improved title summarizing the main substantive takeaway.

    (b) A more informative subtitle describing the sample and variables.

    (c) A caption noting data source, weighting, and key coding decisions.

    (d) At least one direct annotation using `ggrepel` that calls out a key pattern.

    **Answer:** For the final requirement, I transformed the age-turnout analysis into a Lollipop chart to improve the data-to-ink ratio. By removing the bulk of the bars and focusing on the points, the visual noise is significantly reduced. I applied my custom theme, `my_custom_theme`, which follows CRAP principles by centering the title for alignment and using bold text for contrast. Finally, I used `ggrepel` for direct annotation to highlight the 65+ group's peak participation rate.

```
1 # (1) Custom Theme: Building my own reusable style based on CRAP
      principles
2 # I chose theme_minimal as a base to avoid the default gray
      background which can be distracting[cite: 1540].
3 my_custom_theme <- function() {
4   theme_minimal() +
5     theme(
```

```
 6         # CONTRAST & ALIGNMENT: Making the title bold and centered for
      better visual hierarchy [cite: 1526, 1600].
 7         plot.title = element_text(face = "bold", size = 14, hjust =
      0.5),
 8
 9         # HIERARCHY: Using gray text for subtitles to distinguish them
      from the main title [cite: 1854].
10         plot.subtitle = element_text(size = 10, color = "gray30",
      margin = margin(b = 10)),
11         plot.caption = element_text(size = 8, face = "italic", hjust =
      0),
12
13         # REDUCING NON-DATA INK: Removing minor grid lines to focus
      attention on the actual data [cite: 1661, 1741].
14         panel.grid.minor = element_blank(),
15
16         # PROXIMITY: Moving the legend to the bottom for a more
      balanced layout and better spatial alignment [cite: 1529, 1728].
17         legend.position = "bottom"
18       )
19 }
20
21 # (2) Data Prep: Ensuring Logical Factor Ordering
22 # R defaults to alphabetical sorting, so I manually reordered
        categories to follow a logical life-cycle flow.
23 final_plot_data <- ces2015 %>%
24   mutate(age_group = factor(age_group, levels = c("<30", "30-44", "
      45-64", "65+"))) %>%
25   group_by(age_group) %>%
26   summarise(turnout_rate = mean(voted_binary, na.rm = TRUE))
27
28 # (3) Final Visualization: Lollipop Chart Construction
29 # I opted for a Lollipop style because it significantly reduces ink
        compared to bulky bar charts [cite: 342, 343].
30 final_plot <- ggplot(final_plot_data, aes(x = age_group, y = turnout_
      rate)) +
31
32   # ANATOMY: Using points and segments to maintain a high data-to-ink
        ratio [cite: 278, 342].
33   geom_segment(aes(x = age_group, xend = age_group, y = 0, yend =
      turnout_rate), color = "darkblue") +
34   geom_point(size = 5, color = "darkblue") +
35
36   # ANNOTATION: Adding clear labs() for transparency and context [cite
      : 1854, 2165].
37   labs(
38     title = "Electoral Participation Strongly Increases with Age in
      Canada",
39     subtitle = "Analysis based on 2015 CES high-quality respondents",
40     caption = "Source: 2015 Canadian Election Study. Coding reflects
      binary turnout status.",
```

```r
41       x = "Age  Category",
42       y = "Average  Turnout  Probability"
43     ) +
44
45     # DIRECT ANNOTATION: Using  ggrepel  to  highlight  the  most  important
           data  point  without  overlapping[cite:  1852,  1909].
46     geom_text_repel(
47       aes(label  =  ifelse(age_group  ==  "65+",  "Highest  Participation
           Rate",  "")),
48       nudge_y  =  0.05,
49       fontface  =  "bold",
50       segment.color  =  'grey50'
51     ) +
52
53     # Final  Polishing:  Formatting  Y-axis  to  percentage  and  applying  my
           custom  theme[cite:  279,  336].
54     scale_y_continuous(labels  =  scales::percent,  limits  =  c(0,  1)) +
55     my_custom_theme()
56
57 # (4) PDF  Export  for  Report  Integration
58 # Exporting  as  PDF  to  ensure  high-quality  vector  graphics  for  the
       final  LaTeX  document.
59 ggsave("Final_Plot.pdf",  plot  =  final_plot,  width  =  8,  height  =  6)
```



**Electoral Participation Strongly Increases with Age in Canada**

Analysis based on 2015 CES high–quality respondents

*Source: 2015 Canadian Election Study. Coding reflects binary turnout status.*