

Problem Set 1

Data Visualisation for Social Scientists

Due: January 28, 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 January 28, 2026. No late assignments will be accepted.

Roll Call Votes in the European Parliament

Data Manipulation

First, you need to download data from the first six elected European Parliaments on each MEP and how they voted in each recorded roll-call vote.

1. Load these datasets into your global environment:

- **mep_info_26Jul11.xls** (MEP characteristics, EP1–EP5)
- **rcv_ep1.txt** (EP1 roll-call votes)

```
1 install.packages(c("readxl", "readr"))
2 library(readxl)
3 library(readr)
4 library(tidyverse)
5
6 setwd("~/Desktop/whiz1/DataViz_2026/problemSets/PS01/my_answers")
7 mep_info <- read_excel("mep_info_26Jul11.xls", sheet = "EP1", col_names =
  TRUE)
```

```

8 head(mep_info)
9 str(mep_info)
10 rcv_ep1 <- read_table2("rcv_ep1.txt")
11 head(rcv_ep1)
12 str(rcv_ep1)

```

2. Briefly describe (2–3 sentences each) the unit of analysis and key variables in each of these two datasets.

(a) **mep_info.26Jul11.xls (MEP characteristics)**

Unit of analysis: Each row represents an individual Member of the European Parliament (MEP) across the first five European Parliaments (EP1–EP5).

Key variables: The dataset includes MEP identifiers (e.g., `MEP_ID`), personal characteristics (such as name), political affiliation (national party and European Parliament group), country of representation, and ideological positioning variables. This dataset describes attributes of MEPs rather than their voting behavior.

(b) **rcv_ep1.txt (EP1 roll-call votes)**

Unit of analysis: Each row represents a single vote cast by an individual MEP in a recorded roll-call vote during the first European Parliament (EP1).

Key variables: Key variables include the MEP identifier, roll-call vote identifier, and the recorded vote decision (e.g., Yes, No, Abstain, Present but did not vote, or Absent). This dataset captures individual voting behavior on legislative decisions.

3. The `rcv_ep1` data are in a wide format, with `V1`, `V2`, ..., `Vn` as separate vote columns.

- Identify which columns are ID/metadata (*MEPID*, *MEPNAME*, *MS*, *NP*, *EPG*) and which columns are vote decisions ($V_1 \dots V_n$). Tidy the voting data such that each row/observation is a single vote for a single MEP.

```

1 #Question 3
2 rcv_ep1 <- read_delim("rcv_ep1.txt")
3 colnames(rcv_ep1)
4 rcv_ep1_long <- rcv_ep1 %>%
5   pivot_longer(
6     cols = starts_with("V"),
7     names_to = "vote_id",
8     values_to = "decision"
9   )

```

- Create a summary table of counts of decision categories (e.g. Yes/No/Abstain/Present but did not vote/Absent) across all votes.

```

1 rcv_ep1_long <- rcv_ep1_long %>%
2   mutate(
3     decision_label = case_when(
4       decision == 1 ~ "Yes",
5       decision == 2 ~ "No",
6       decision == 3 ~ "Abstain",
7       decision == 4 ~ "Absent",
8       decision == 5 ~ "Present but did not vote",
9       decision == 0 ~ "Not a member",
10      TRUE ~ "Other"
11    )
12  )
13
14 decision_summary <- rcv_ep1_long %>%
15   count(decision_label) %>%
16   arrange(desc(n))
17
18 decision_summary

```

Decision	Count
Absent	109,224
Present but did not vote	103,618
Not a member	99,753
Yes	88,185
No	75,171
Abstain	9,577

Table 1: Distribution of vote decisions in EP1 roll-call votes

4. Construct a new dataset that combines MEP-level information with their vote decisions from EP1 in long format (from part 3). Check for missingness.

```

1 #Question 4
2 meta_cols <- c("MEPID", "MEPNAME", "MS", "NP", "EPG")
3 ep1_combined <- rcv_ep1_long %>%
4   select(all_of(meta_cols), vote_id, decision, decision_label)
5 head(ep1_combined)
6 str(ep1_combined)
7 missing_summary <- ep1_combined %>%
8   summarise(across(everything(), ~sum(is.na(.))))
9
10 missing_summary

```

MEPID	MEPNAME	MS	NP	EPG	vote_id	decision	decision_label
0		0	0	0	0	0	0

Table 2: Example observation from the combined EP1 voting dataset

5. Compute, for each EP group in EP1:

- The mean rate of Yes votes (Yes over Yes+No+Abstain) across all roll calls.

```

1 #Question 5.1
2 #Filter only relevant votes (Yes, No, Abstain)
3 votes_filtered <- ep1_combined %>%
4   filter(decision_label %in% c("Yes", "No", "Abstain"))
5
6 #Compute Yes rate per roll call per EP group
7 yes_rate_per_vote <- votes_filtered %>%
8   group_by(EPG, vote_id) %>% # each vote within each EP group
9   summarise(
10     yes_count = sum(decision_label == "Yes"),
11     total_votes = n(),
12     yes_rate = yes_count / total_votes,
13     .groups = "drop"
14   )
15
16 #Compute mean Yes rate across all roll calls for each EP group
17 mean_yes_rate_epg <- yes_rate_per_vote %>%
18   group_by(EPG) %>%
19   summarise(mean_yes_rate = mean(yes_rate), .groups = "drop") %>%
20   arrange(desc(mean_yes_rate))
21
22 mean_yes_rate_epg

```

EP Group	Mean Yes Rate
N	0.613
S	0.563
E	0.532
G	0.514
M	0.512
L	0.495
R	0.486
C	0.437

Table 3: Average proportion of Yes votes by EP group in EP1

- The mean abstention rate.

```

1 #Question5.2
2 #Filter only relevant votes (Yes, No, Abstain)
3 votes_filtered <- ep1_combined %>%
4   filter(decision_label %in% c("Yes", "No", "Abstain"))
5
6 #Compute Abstain rate per roll call per EP group
7 abstain_rate_per_vote <- votes_filtered %>%
8   group_by(EPG, vote_id) %>%
9   summarise(
10     abstain_count = sum(decision_label == "Abstain"),
11     total_votes = n(),
12     abstain_rate = abstain_count / total_votes,
13     .groups = "drop"
14   )
15
16 #Compute mean Abstain rate across all roll calls for each EP group
17 mean_abstain_rate_epg <- abstain_rate_per_vote %>%
18   group_by(EPG) %>%
19   summarise(mean_abstain_rate = mean(abstain_rate), .groups = "drop")
20   %>%
21   arrange(desc(mean_abstain_rate))
22
23 # View the result
24 mean_abstain_rate_epg

```

EP Group	Mean Abstention Rate
R	0.2690
M	0.0893
G	0.0806
C	0.0725
L	0.0654
S	0.0598
N	0.0580
E	0.0215

Table 4: Average proportion of abstentions by EP group in EP1

- The mean vote preferences along the two contested dimensions (NOM-D1 and NOM-D2).

```

1 #Question5.3
2 library(dplyr)
3 mep_info <- read_excel("mep_info_26Jul11.xls", sheet = "EP1",
4   col_names = TRUE)
5
6 #clean column names
7 colnames(mep_info) <- colnames(mep_info) %>%
8   trimws() %>%
9   gsub(" ", ".", .) # replace spaces with dots

```

```

9
10 #Ensure the key columns exist
11 colnames(mep_info)
12 head(mep_info)
13
14
15 #Fix column types and trim spaces
16 mep_info <- mep_info %>%
17   mutate(
18     MEP.id = trimws(as.character(MEP.id)), # ensure same type as
19     ep1_combined$MEPID
20     'NOM-D1' = as.numeric('NOM-D1'),
21     'NOM-D2' = as.numeric('NOM-D2')
22   )
23 ep1_combined <- ep1_combined %>%
24   mutate(MEPID = trimws(as.character(MEPID)))
25
26
27 #Merge NOM-D1 and NOM-D2 into long-format vote dataset
28 ep1_combined_nom <- ep1_combined %>%
29   left_join(
30     mep_info %>% select(MEP.id, 'NOM-D1', 'NOM-D2'),
31     by = c("MEPID" = "MEP.id")
32   )
33
34
35 #Compute mean NOM-D1 and NOM-D2 per EP group
36 mean_nominate_epg <- ep1_combined_nom %>%
37   group_by(EPG) %>%
38   summarise(
39     mean_nom_d1 = mean('NOM-D1', na.rm = TRUE),
40     mean_nom_d2 = mean('NOM-D2', na.rm = TRUE),
41     .groups = "drop"
42   ) %>%
43   arrange(mean_nom_d1)
44
45
46 #View the result
47 mean_nominate_epg

```

EP Group	Mean NOM-D1	Mean NOM-D2
R	-0.586	-0.0419
M	-0.357	-0.2010
S	-0.0980	0.2610
N	0.250	-0.3860
G	0.280	-0.8180
L	0.409	-0.3240
E	0.512	-0.2770
C	0.811	0.5300
0	NaN	NaN

Table 5: Average NOMINATE coordinates (Dimension 1 and Dimension 2) by EP group in EP1

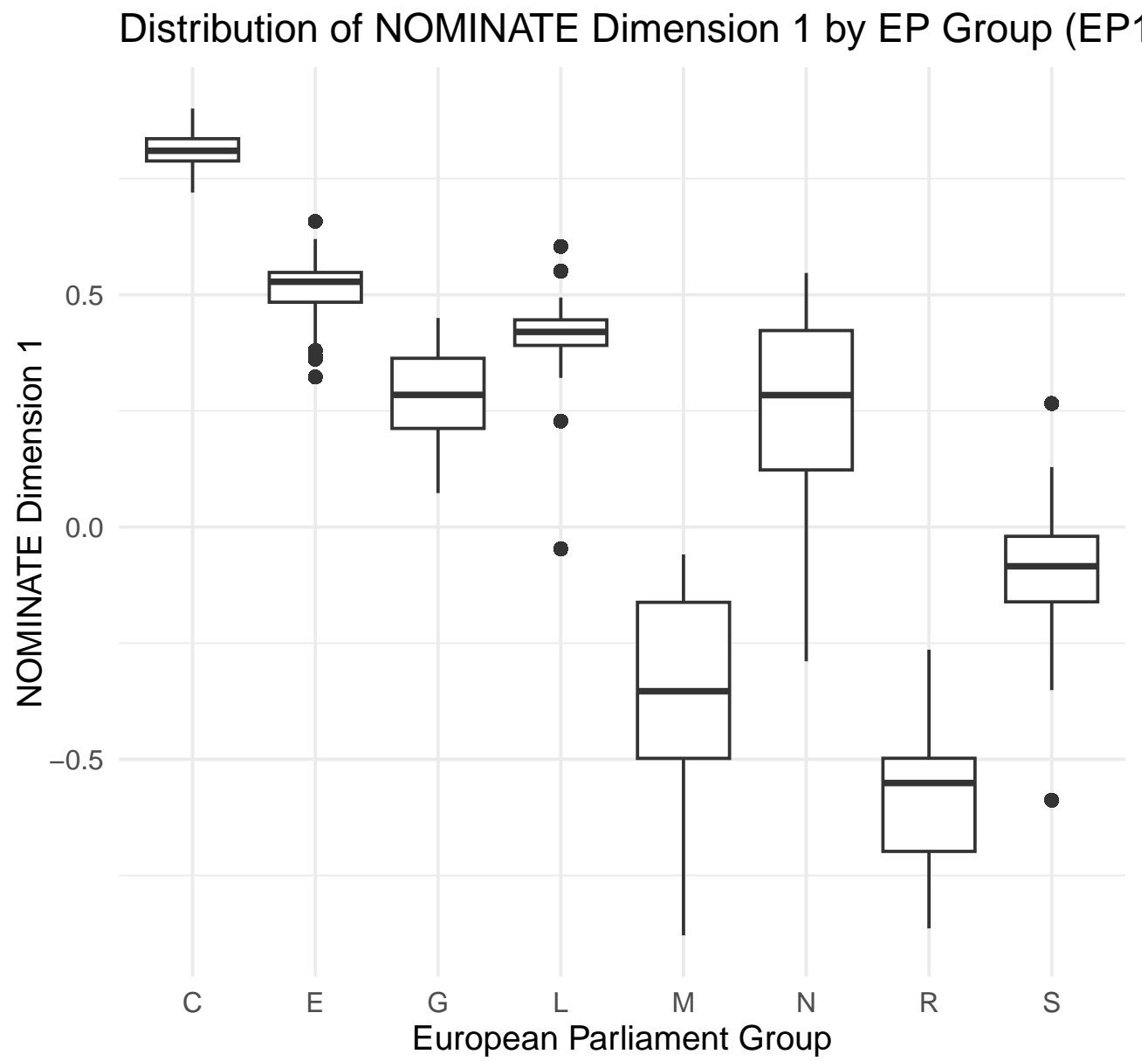
Data Visualization

1. Plot the distribution of the first NOMINATE dimension by EP group, and explain any trends you see.

```

1 #Data Visualization
2 #Question 1
3 library(ggplot2)
4 #Keep only observations with valid NOM-D1 values
5 plot_data <- ep1_combined_nom %>%
6   filter(!is.na('NOM-D1'))
7
8 #Boxplot of NOM-D1 by EP group
9 ggplot(plot_data, aes(x = EPG, y = 'NOM-D1')) +
10   geom_boxplot() +
11   labs(
12     title = "Distribution of NOMINATE Dimension 1 by EP Group (EP1)",
13     x = "European Parliament Group",
14     y = "NOMINATE Dimension 1"
15   ) +
16   theme_minimal()

```



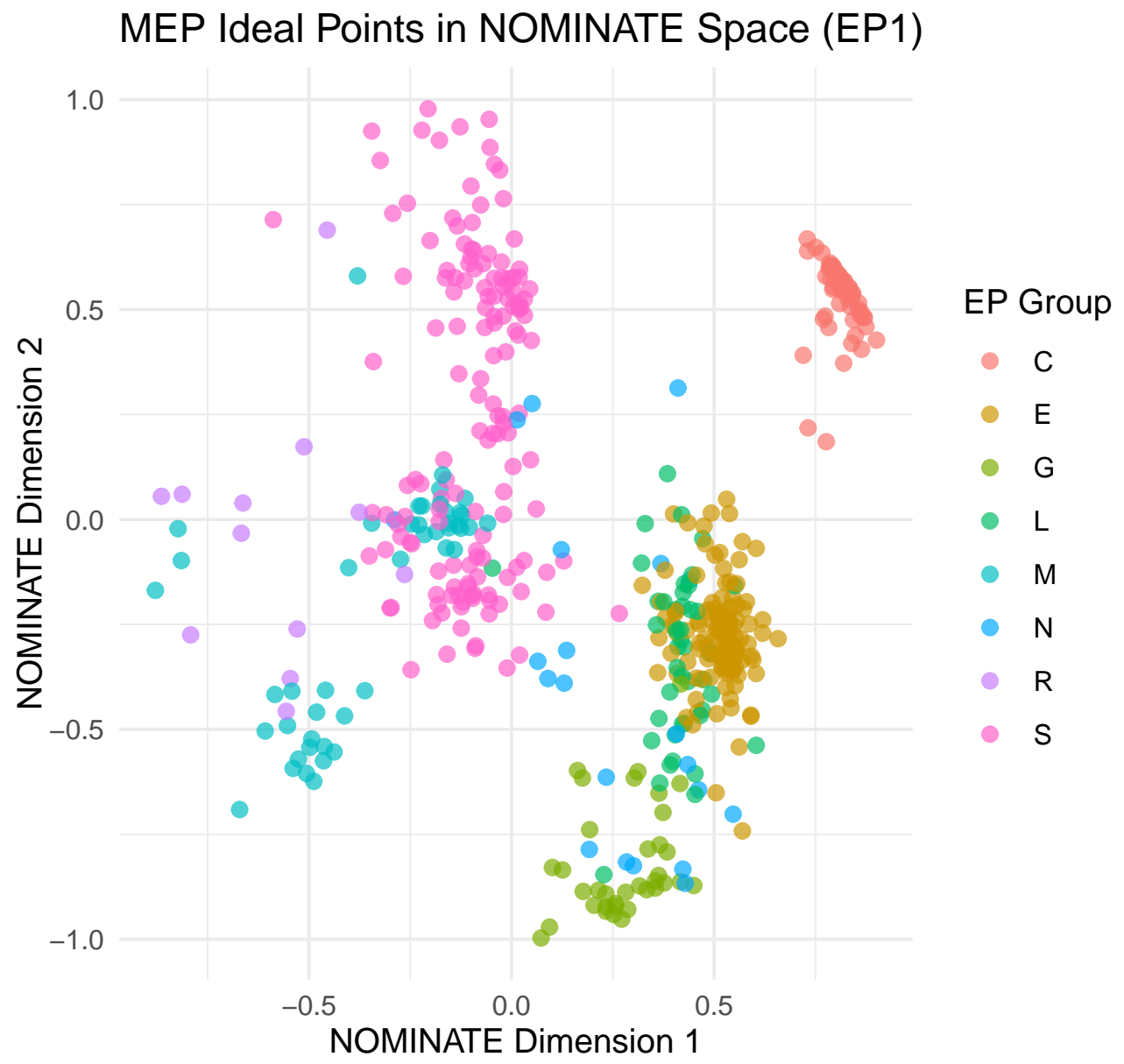
The plot shows clear differences in NOM-D1 across EP groups, with left-leaning groups concentrated at lower values and right-leaning groups at higher values. There is a negative linear relationship shown in the graph. This indicates that the first NOMINATE dimension captures an underlying ideological cleavage in the European Parliament.

2. Make a scatterplot of *nomdim1* (x-axis) and *nomdim2* (y-axis), with one point per MEP and color by EP group.

```

1 #Question 2
2 #Create a unique MEP-level dataset
3 mep_level <- ep1_combined_nom %>%
4   select(MEPID, EPG, 'NOM-D1', 'NOM-D2') %>%
5   distinct() %>%                                # one row per MEP
6   filter(!is.na('NOM-D1'), !is.na('NOM-D2'))
7
8 #Scatterplot
9 ggplot(mep_level, aes(x = 'NOM-D1', y = 'NOM-D2', color = EPG)) +
10   geom_point(alpha = 0.7, size = 2) +
11   labs(
12     title = "MEP Ideal Points in NOMINATE Space (EP1)",
13     x = "NOMINATE Dimension 1",
14     y = "NOMINATE Dimension 2",
15     color = "EP Group"
16   ) +
17   theme_minimal()

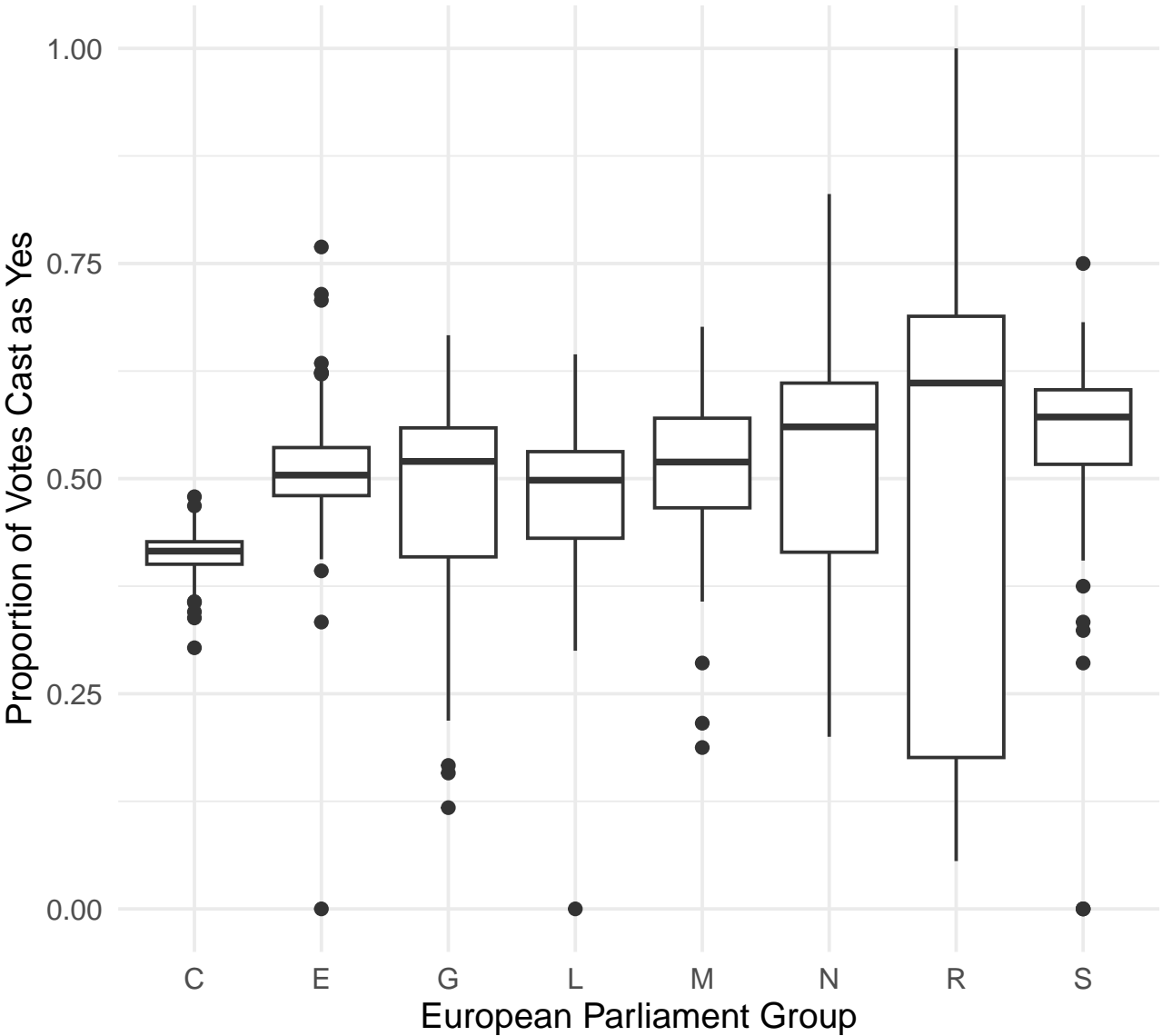
```



3. Produce a boxplot of the proportion voting *Yes* by EP group to visualize cohesion.

```
1 #Question 3
2 # Step 1: Compute proportion of Yes votes per MEP
3 mep_yes_rate <- epl_combined %>%
4   filter(decision_label %in% c("Yes", "No", "Abstain")) %>% #
5     denominator
6   group_by(MEPID, EPG) %>%
7   summarise(
8     yes_rate = mean(decision_label == "Yes"),
9     .groups = "drop"
10  )
11 #Boxplot by EP group
12 ggplot(mep_yes_rate, aes(x = EPG, y = yes_rate)) +
13   geom_boxplot() +
14   labs(
15     title = "Distribution of Yes Vote Proportions by EP Group (EP1)",
16     x = "European Parliament Group",
17     y = "Proportion of Votes Cast as Yes"
18   ) +
19   theme_minimal()
```

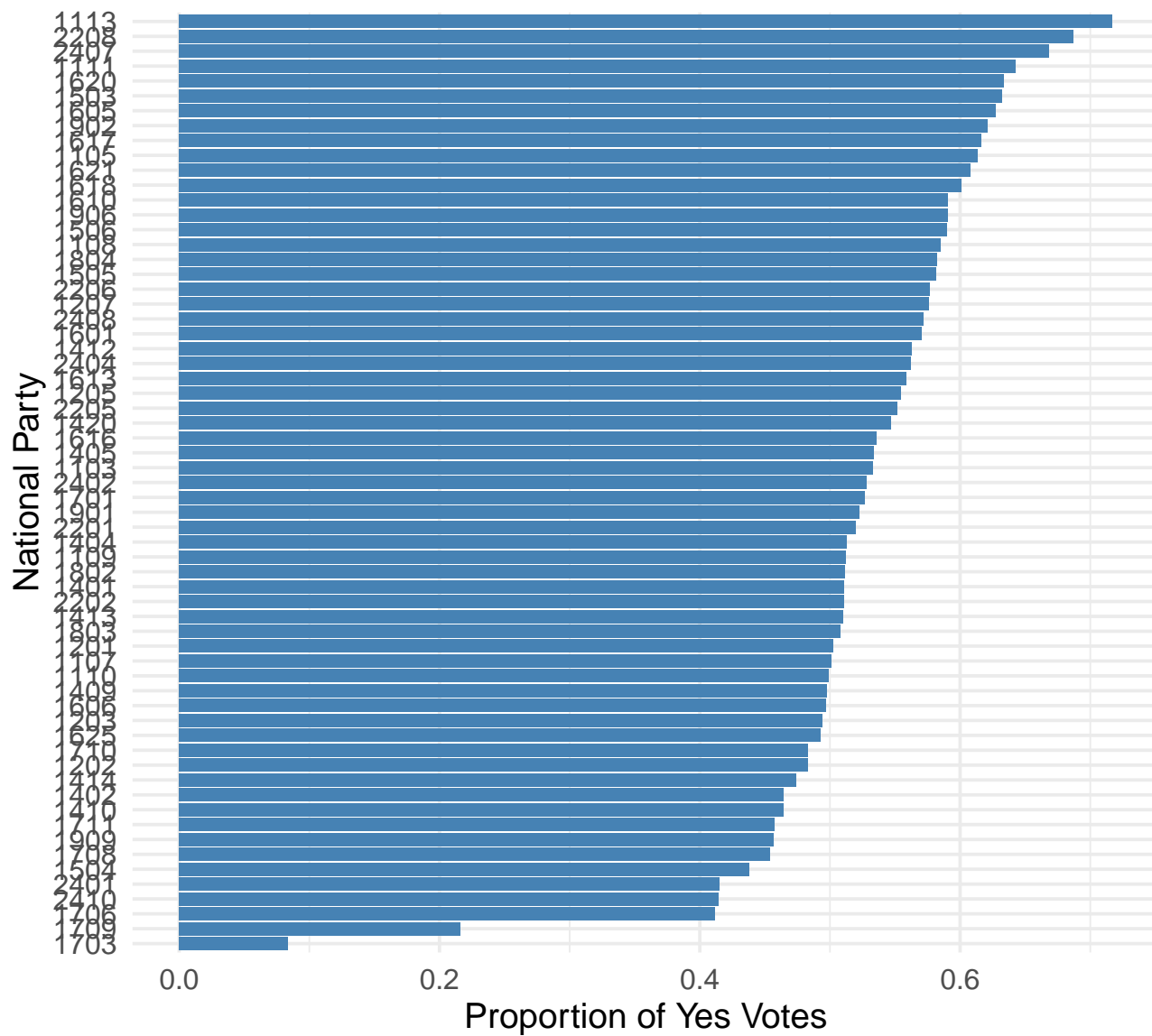
Distribution of Yes Vote Proportions by EP Group (EP1)



4. Display the proportion voting *Yes* per year by national party using a bar plot.

```
1 #Question4
2 ep1_combined %>%
3   filter(decision_label %in% c("Yes", "No", "Abstain")) %>%
4   group_by(NP) %>%
5   summarise(
6     yes_rate = mean(decision_label == "Yes"),
7     .groups = "drop"
8   ) %>%
9   ggplot(aes(x = reorder(factor(NP), yes_rate), y = yes_rate)) +
10  geom_col(fill = "steelblue") +
11  coord_flip() +
12  labs(
13    title = "Proportion Voting Yes by National Party (EP1)",
14    x = "National Party",
15    y = "Proportion of Yes Votes"
16  ) +
17  theme_minimal()
```

Proportion Voting Yes by National Party (EP1)



5. For each EP group, calculate the average *Yes* share per year and plot a line graph.