

GROUP PRESENTATION

Andrea Baños, Alejandra Costa, Eric Hausken, María Jurado

SPANISH ELECTIONS ANALYSIS

This project will analyze **electoral and survey data**, from 2008 to 2019, from the Spanish Congress of Deputies.

The following R packages will be in our report:

Libraries:

```
1 library(tidyverse)
2 library(lubridate)
3 library(glue)
4 library(dplyr)
5 library(ggplot2)
6 library(corrplot)
7 library(forcats)
8 library(patchwork)
9 library(ghibli)
```

1. Tidy data

1.1 Tidy data: Election data

We first prepare and clean election dataset by applying some transformations.

```
1 election_pivot <- election_data |>
2   pivot_longer(
3     cols = `BERDEAK-LOS VERDES` : `COALICIÓN POR MELILLA` ,
4     names_to = "party",
5     values_to = "votos"
6   ) |>
7   drop_na(votos) |>
8   select(-c(vuelta, tipo_eleccion, codigo_distrito_electoral)) |>
9   mutate(
10     date_elec_ym = lubridate::ym(paste(anno, mes)),
11     .before = anno
12   ) |>
13   mutate(
14     codigo_ccaa = as_factor(codigo_ccaa),
15     codigo_municipio = as_factor(codigo_municipio),
16     codigo_provincia = as_factor(codigo_provincia),
17     municipio = as_factor(paste(codigo_ccaa,
18                                 codigo_provincia, codigo_municipio,
19                                 sep = "-")),
20     party = as_factor(party)
21   ) |>
22   select(-c(mes, anno, numero_mesas, participacion_1, participacion_2))
```

1.1 Tidy data: Election data

We also grouped the parties into the main classifications.

```
1 election_pivot <- election_pivot |>
2
3 mutate(party = case_when(
4
5     str_detect(party, "PODEMOS") |
6     str_detect(party, "PODEM") |
7     str_detect(party, "VERDES") |
8     str_detect(party, "IZQUIERDA UNIDA") |
9     str_detect(party, "ESQUERRA UNIDA") |
10    str_detect(party, "EZKER BATUA") ~ "PODEMOS",
11
12    str_detect(party, "SOCIALISTA") |
13    str_detect(party, "SOCIALISTES") ~ "PARTIDO SOCIALISTA OBRERO ESPAÑOL",
14
15    str_detect(party, "PARTIDO POPULAR") ~ "PARTIDO POPULAR",
16
17    str_detect(party, "CIUDADANÍA") |
18    str_detect(party, "CIUDADANIA") ~ "CIUDADANOS-PARTIDO DE LA CIUDADANIA",
19
20    str_detect(party, "NACIONALISTA VASCO") ~ "EUZKO ALDERDI JELTZALEA-PARTIDO NACIONALISTA VASCO",
21
22    str_detect(party, "NACIONALISTA GALEGO") ~ "BLOQUE NACIONALISTA GALEGO",
23
24    str_detect(party, "MÉS COMPROMÍS") ~ "MÉS COMPROMÍS",
25
```

1.1 Tidy data: Election data

After cleaning the dataset and selecting just the relevant information, the election data look as follows:

```
1 head(election_pivot)
```

```
# A tibble: 6 × 11
  date_elec_ym codigo_ccaa codigo_provincia codigo_municipio censo votos_blanco
  <date>      <fct>      <fct>      <fct>      <dbl>      <dbl>
1 2008-03-01   14         01         001         1838         23
2 2008-03-01   14         01         001         1838         23
3 2008-03-01   14         01         001         1838         23
4 2008-03-01   14         01         001         1838         23
5 2008-03-01   14         01         001         1838         23
6 2008-03-01   14         01         001         1838         23
# i 5 more variables: votos_nulos <dbl>, votos_candidaturas <dbl>, party <chr>,
#   votos <dbl>, municipio <fct>
```

1.2 Tidy data: Abbreviation data

The abbreviations are unified in the `abbrev` table, and those corresponding to non-relevant parties are categorized as “OTHERS”

As a result, we get a reference table containing unique parties, with their corresponding abbreviation.

```
1 print(abbrev)
```

```
# A tibble: 13 × 2
  denominacion                siglas
  <chr>                  <chr>
1 OTHERS                OTHERS
2 EUZKO ALDERDI JELTZALEA-PARTIDO NACIONALISTA VASCO EAJ-PNV
3 PODEMOS                PODEMOS
4 PARTIDO SOCIALISTA OBRERO ESPAÑOL    PSOE
5 PARTIDO POPULAR        PP
6 ESQUERRA REPUBLICANA DE CATALUNYA    ERC
7 CONVERGENCIA I UNIO    CiU
8 BLOQUE NACIONALISTA GALEGO    BNG
9 CIUDADANOS-PARTIDO DE LA CIUDADANIA  CS
10 EUSKAL HERRIA BILDU    EH-BILDU
11 VOX                    VOX
12 MÉS COMPROMÍS        COMPROMIS
13 MÁS PAÍS              M PAÍS
```

1.2 Tidy data: Abbreviation data

A vector is created with the final abbreviations, which will be useful later when cleaning and preparing the survey data.

```
[1] "OTHERS"      "EAJ-PNV"      "PODEMOS"      "PSOE"         "PP"           "ERC"  
[7] "CiU"         "BNG"          "CS"           "EH-BILDU"     "VOX"          "COMPROMIS"  
[13] "M PAÍS"
```

Lastly, we join the abbreviation data with the election data. By doing so, we get the corresponding abbreviations of the parties in the election dataset and get rid of the variable “Party”, containing the whole name.

```
1 election_pivot <- election_pivot |>  
2   left_join(abbrev, by = c("party" = "denominacion"))
```


1.3 Tidy data: Survey data

We apply the following transformations to satisfy the specified conditions:

```
1 surveys <- surveys |>
2   filter(year(date_elec) >= 2008 ) |>
3   filter(exit_poll == FALSE) |>
4   filter(size > 750) |>
5   mutate(fieldwork_days = field_date_to - field_date_from) |>
6   filter(fieldwork_days > 1)
7
8 surveys <- surveys |>
9   mutate(date_elec_ym = ym(paste(lubridate::year(surveys$date_elec),
10                                lubridate::month(surveys$date_elec))),
11         .after = date_elec)
```

Some other transformations are applied, so that we can work with a more organised dataset:

```
1 surveys_pivot <- surveys |>
2   pivot_longer(cols = UCD:EV,
3               names_to = "party",
4               values_to = "votes_percent") |>
5   drop_na(votes_percent) ## removed rows with zero (NA) votes for that party
```

1.3 Tidy data: Survey data

Survey data is aligned with the rest of the datasets.

```
1 surveys_pivot<- surveys_pivot |>
2   mutate(party = ifelse(
3     party %in% vector_abbrev, party, "OTHER")) |>
4   select(-c(type_survey,exit_poll, id_pollster, media))
```

2. Questions

- How is the vote of national parties (PSOE, PP, VOX, CS, MP, UP - IU) distributed against regional or nationalist parties?
- Which party was the winner in the municipalities with more than 100,000 habitants (census) in each of the elections?
- Which party was the second when the first was the PSOE? And when the first was the PP?
- Who benefits from low turnout?
- How to analyze the relationship between census and vote? Is it true that certain parties win in rural areas
- How to calibrate the error of the polls (remember that the polls are voting intentions at national level)?
- In which election were the polls most wrong?
- How were the polls wrong in national parties (PSOE, PP, VOX, CS, MP, UP - IU)?
- Which polling houses got it right the most and which ones deviated the most from the results?

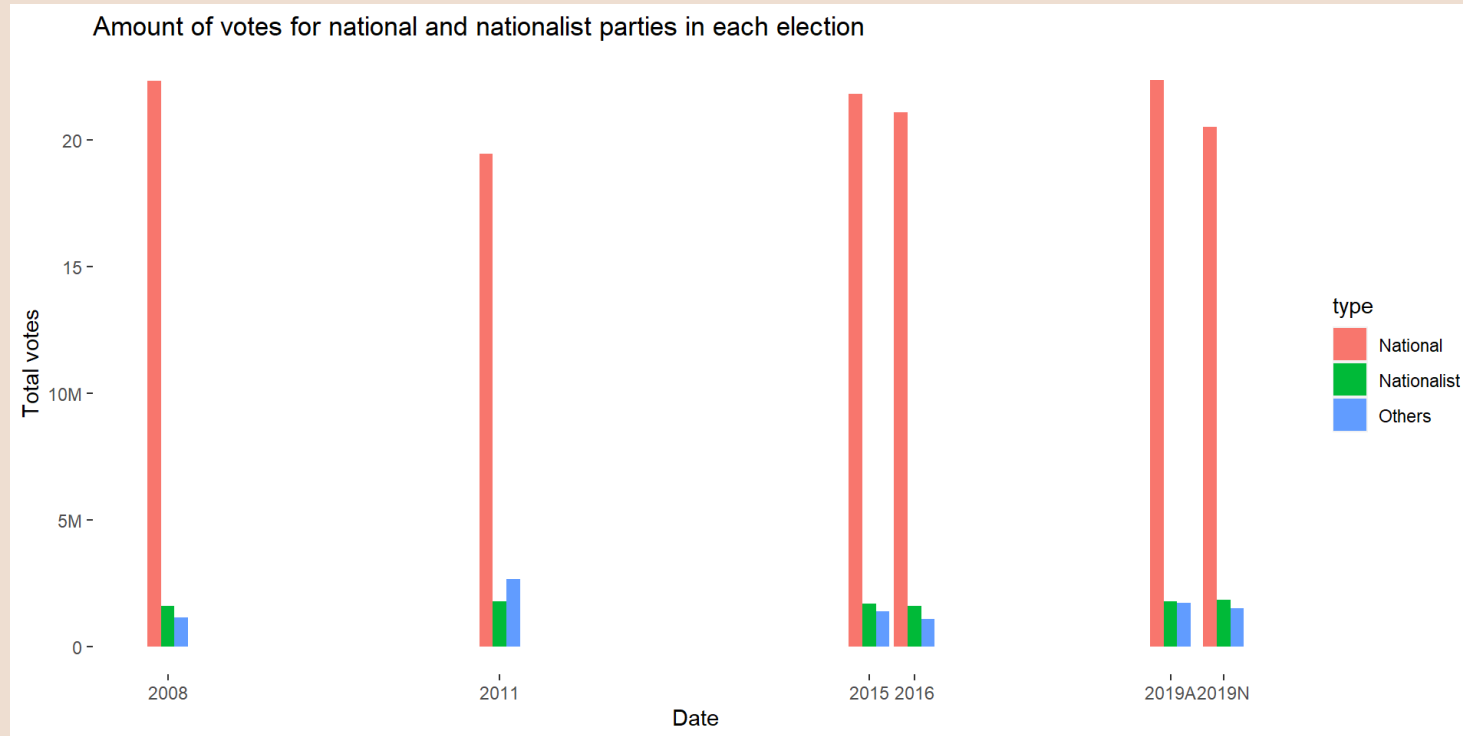
2.1 How is the vote of national parties distributed against regional or nationalist parties?

2.1 How is the vote of national parties distributed against regional or nationalist parties?

- First step is to group the parties into national and nationalists parties.
- NATIONAL: PSOE, PP, VOX, CS, Más País, Podemos
- NATIONALIST: PNV, Bloque Nacionalista Galego, Mès Compromís, Covergencia i Unió, Esquerra Republicana de Catalunya and Bildu.
- Total votes by type

```
1 data_new<-  
2   election_pivot |>  
3   group_by(type,date_elec_ym) |>  
4   summarise(total_votes = sum(votos, na.rm = TRUE)) |>  
5   ungroup()
```

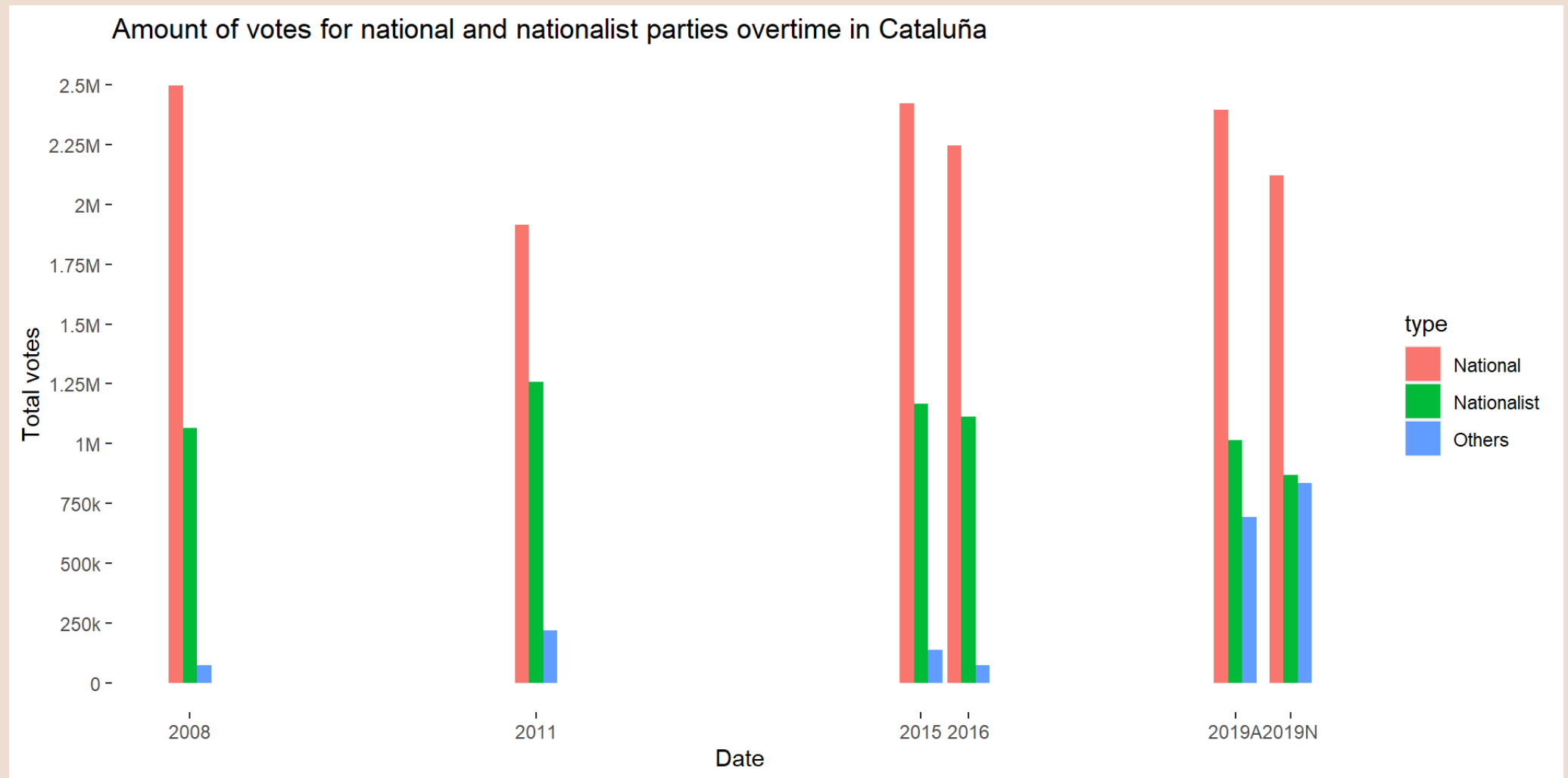
2.1 How is the vote of national parties distributed against regional or nationalist parties?



This graph is really informative.

- Total participation in each election.
- Evolution in votes for all 3 types of parties
- Compare among different elections.

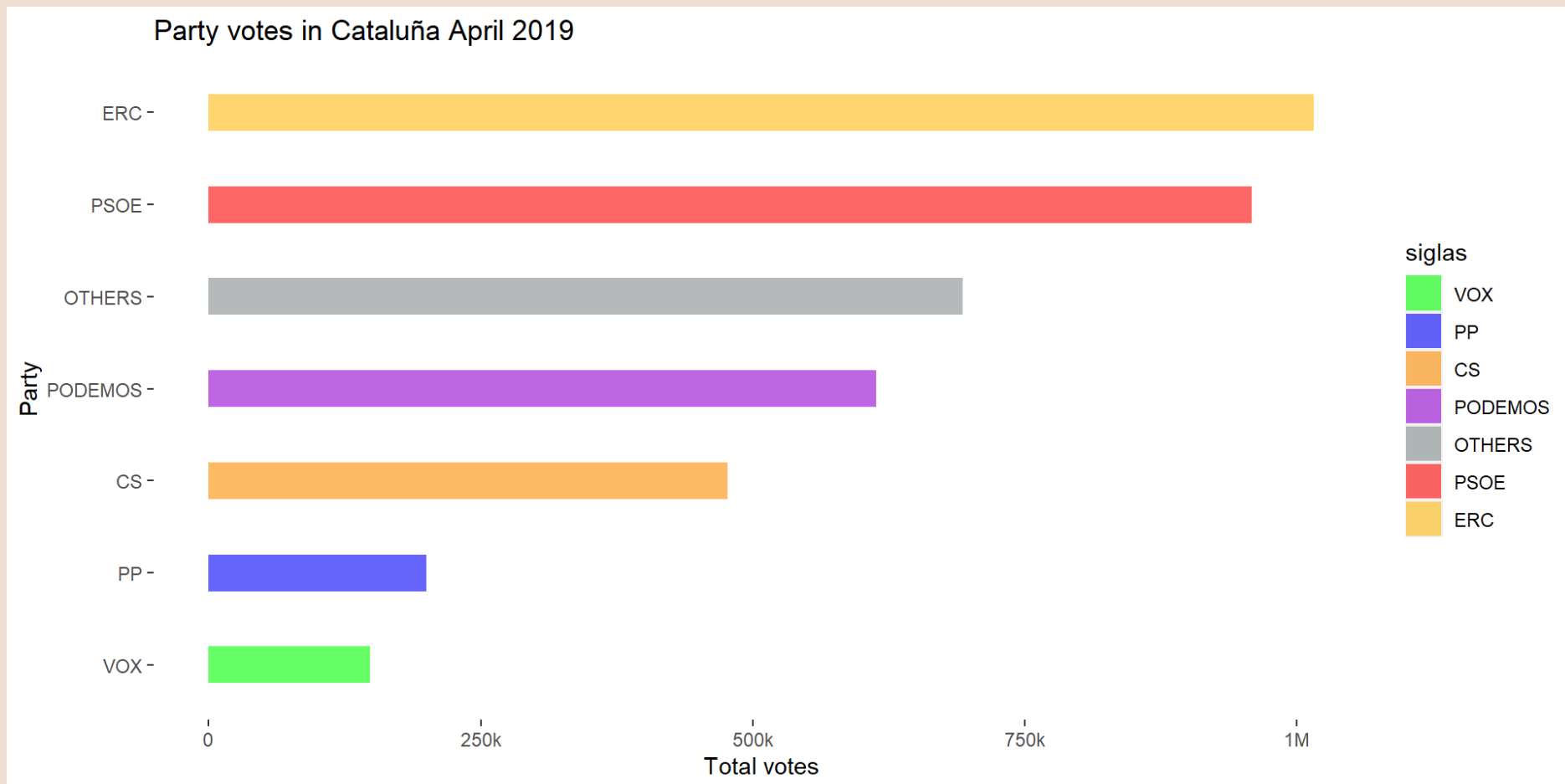
In Catalunya



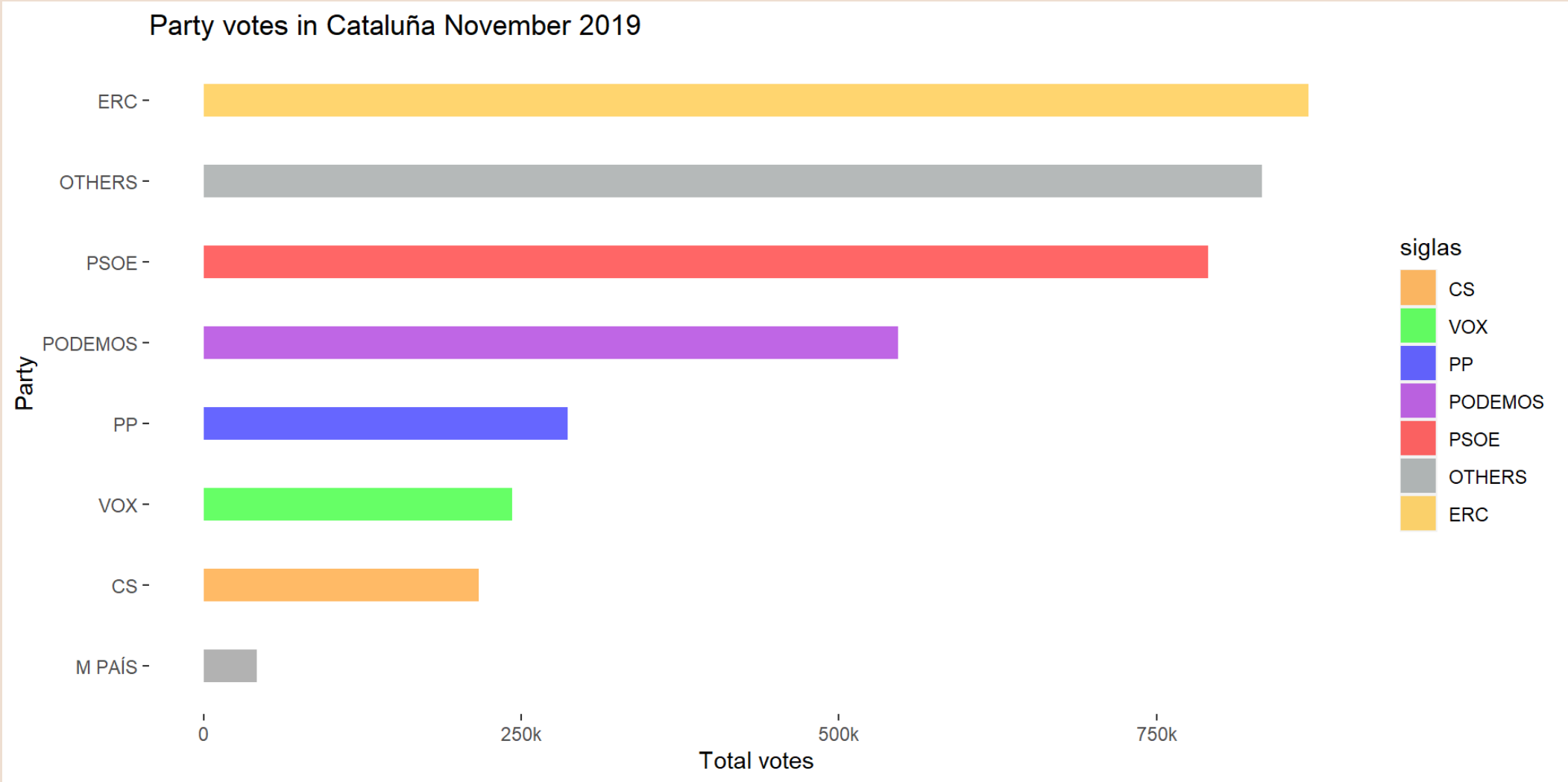
- Decrease in nationalism.
- Bias: Junts is not considered as nationalist → could explain the increase in other parties.

Evolution of party votes in Cataluña in April 2019

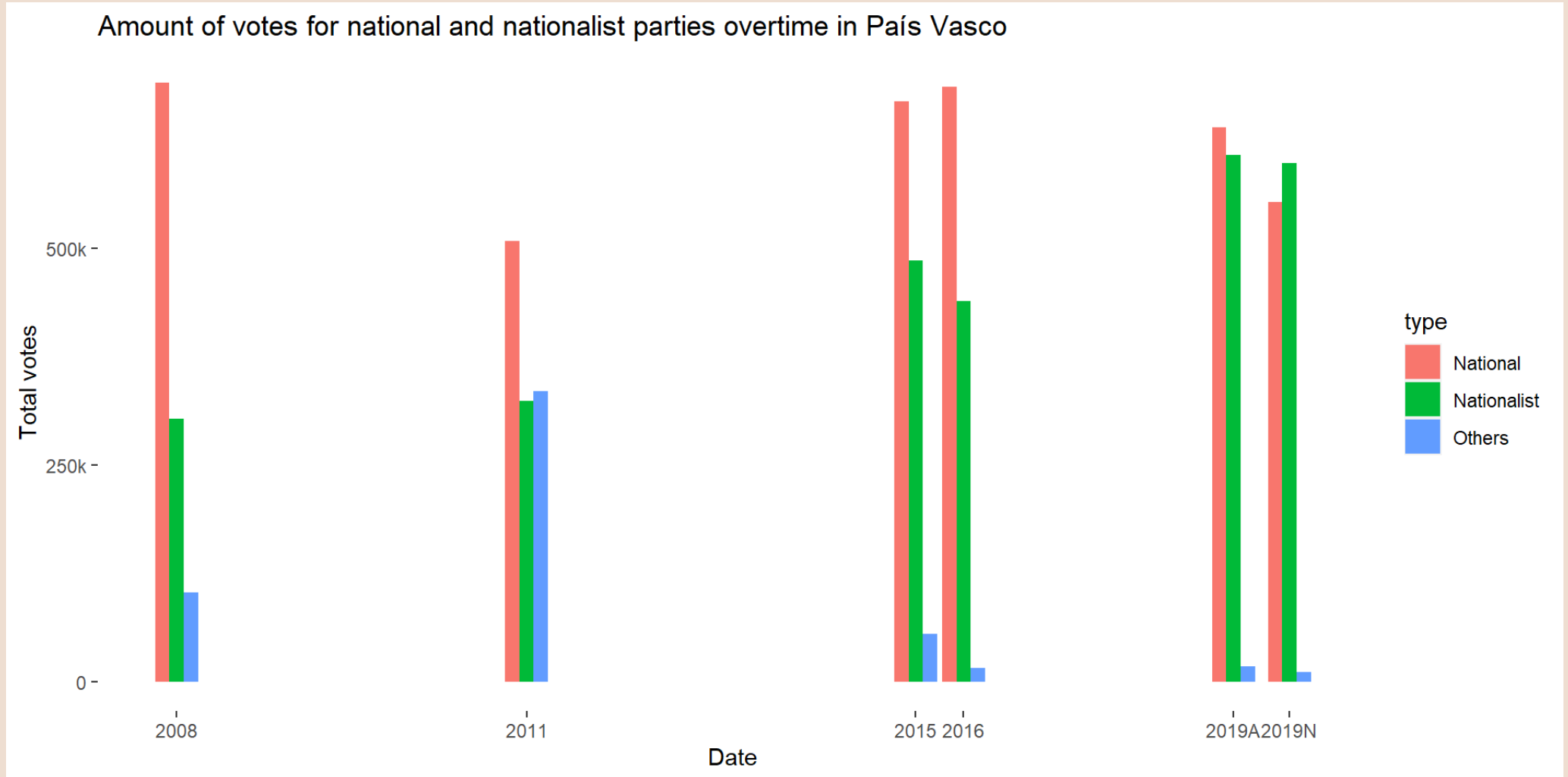
- Using library(forcats)
- Order from most voted to less voted party



Evolution of party votes in Cataluña in November 2019

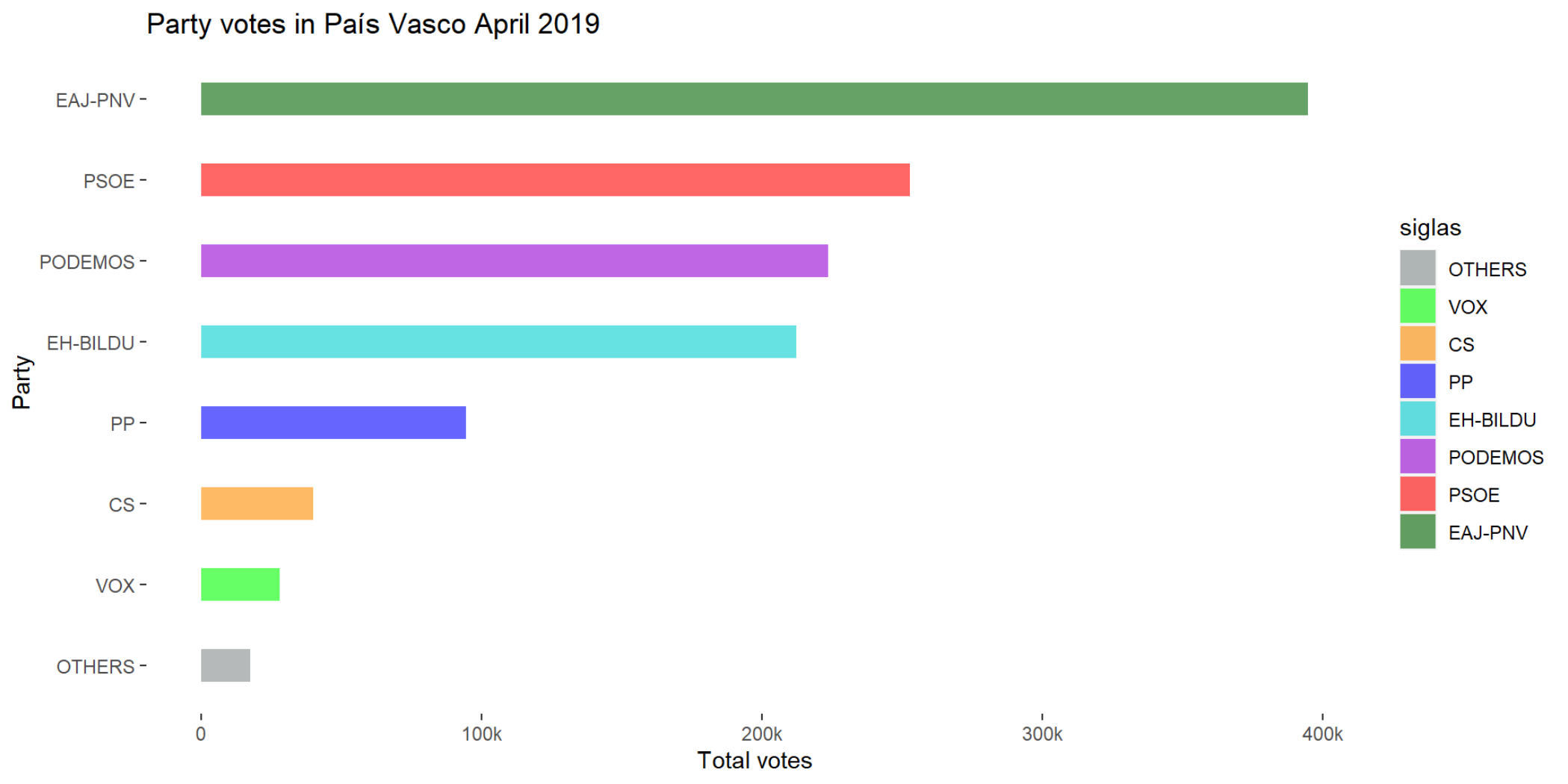


In País Vasco

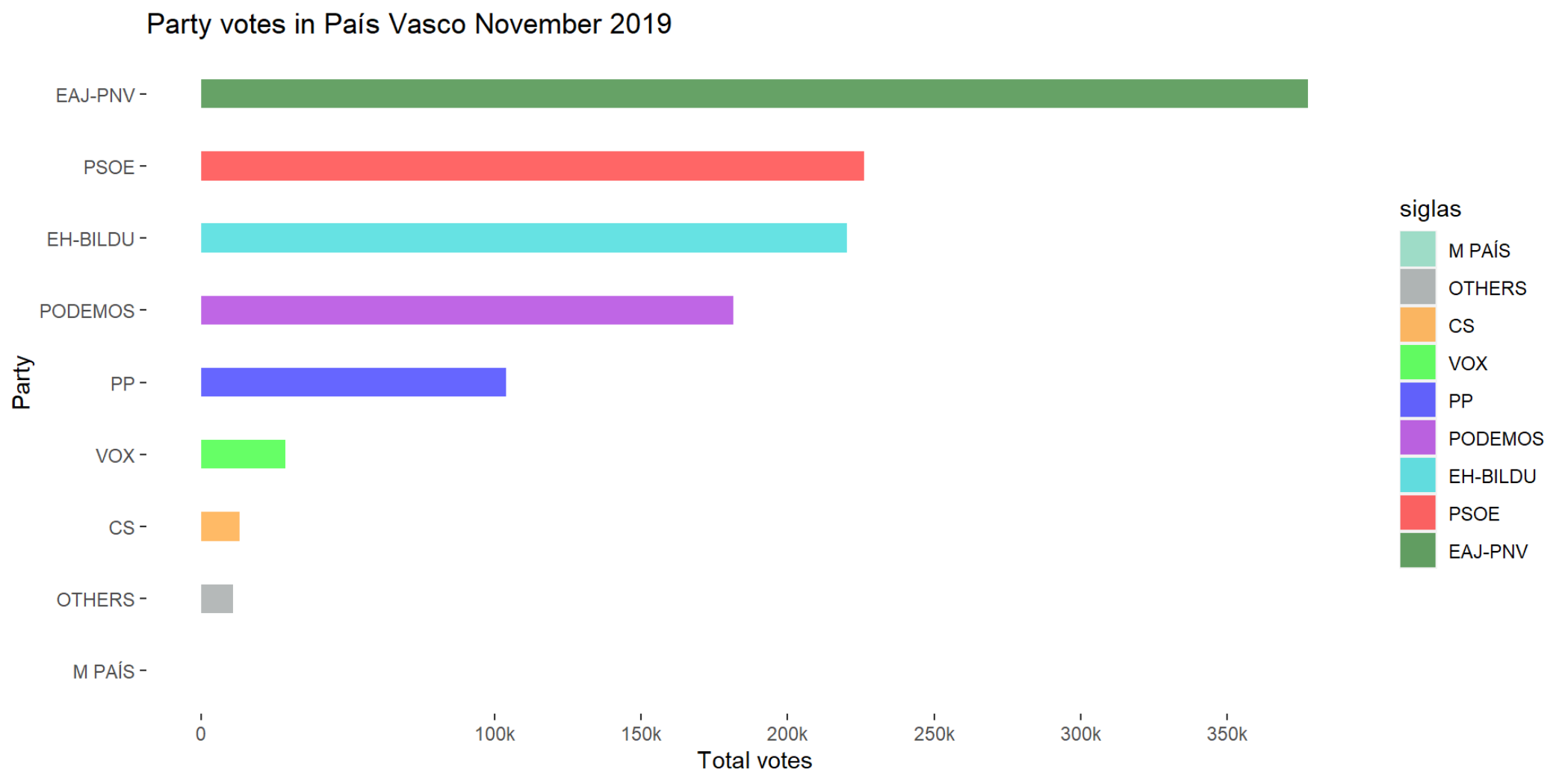


- Increase in nationalism overtime.
- Other parties have decreased significantly the amount of votes.

Party votes País Vasco April 2019



Party votes País Vasco November 2019



2.2 Which party was the winner in the municipalities with more than 100,000 habitants (census) in each of the elections?

2.2 Which party was the winner in the municipalities with more than 100,000 habitants (census) in each of the elections?

```
1 winners <- election_pivot |>
2   filter(censo>100000) |>
3   group_by(codigo_municipio, date_elec_ym, party) |>
4   summarise(total_votes = sum(votos, na.rm = TRUE)) |>
5   slice(which.max(total_votes)) |>
6   ungroup()
7
8 winners
```

A tibble: 231 × 4

| | codigo_municipio <fct> | date_elec_ym <date> | party <chr> | total_votes <dbl> |
|----|---------------------------|------------------------|-----------------------------------|----------------------|
| 1 | 003 | 2008-03-01 | PARTIDO POPULAR | 49909 |
| 2 | 003 | 2011-11-01 | PARTIDO POPULAR | 55858 |
| 3 | 003 | 2015-12-01 | PARTIDO POPULAR | 36149 |
| 4 | 003 | 2016-06-01 | PARTIDO POPULAR | 38470 |
| 5 | 003 | 2019-04-01 | PARTIDO SOCIALISTA OBRERO ESPAÑOL | 28729 |
| 6 | 003 | 2019-11-01 | PARTIDO SOCIALISTA OBRERO ESPAÑOL | 27074 |
| 7 | 013 | 2008-03-01 | PARTIDO POPULAR | 49463 |
| 8 | 013 | 2011-11-01 | PARTIDO POPULAR | 53152 |
| 9 | 013 | 2015-12-01 | PARTIDO POPULAR | 34111 |
| 10 | 013 | 2016-06-01 | PARTIDO POPULAR | 38809 |

i 221 more rows

**2.3 Which party was the second
when the first was the PSOE?
And when the first was the PP?**

2.3 Which party was the second when the first was the PSOE? And when the first was the PP?

```
1 winners2 <- election_pivot |>
2   group_by(date_elec_ym, party) |>
3   summarise(total_votes = sum(votos, na.rm = TRUE)) |>
4   slice_max(total_votes, n = 2) |>
5   ungroup()
6
7 winners2
```

```
# A tibble: 12 × 3
  date_elec_ym party          total_votes
  <date>      <chr>          <dbl>
1 2008-03-01  PARTIDO SOCIALISTA OBRERO ESPAÑOL 11078605
2 2008-03-01  PARTIDO POPULAR          10171828
3 2011-11-01  PARTIDO POPULAR          10838951
4 2011-11-01  PARTIDO SOCIALISTA OBRERO ESPAÑOL 6987723
5 2015-12-01  PARTIDO POPULAR          7114123
6 2015-12-01  PODEMOS                  5640709
7 2016-06-01  PARTIDO POPULAR          7800328
8 2016-06-01  PARTIDO SOCIALISTA OBRERO ESPAÑOL 5424130
9 2019-04-01  PARTIDO SOCIALISTA OBRERO ESPAÑOL 7481667
10 2019-04-01  PARTIDO POPULAR          4356714
11 2019-11-01  PARTIDO SOCIALISTA OBRERO ESPAÑOL 6752314
12 2019-11-01  PARTIDO POPULAR          5021622
```


2.4 Who benefits from low turnout?

2.4 Who benefits from low turnout?

In order to calculate the low turnout in the surveys, we divided the turnout by its mean to generate a ratio. If this ratio exceeds 1, it indicates high turnout, and if it is less than 1, it suggests low turnout.

```
1 summary(surveys_pivot$turnout)
```

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. | NA's |
|-------|---------|--------|-------|---------|-------|------|
| 59.40 | 66.50 | 70.00 | 69.84 | 72.70 | 79.90 | 8204 |

```
1 summary(surveys_pivot$votes_percent)
```

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|------|---------|--------|-------|---------|-------|
| 0.00 | 1.30 | 4.10 | 10.49 | 18.60 | 49.30 |

```
1 surveys_pivot$low_turnout <- surveys_pivot$turnout/68.38
```

2.4 Who benefits from low turnout?

Therefore, a data frame named `low_turnout` has been created, which is defined as turnout below the mean. The parties that benefit from it are those with more votes than the mean within the low turnout.

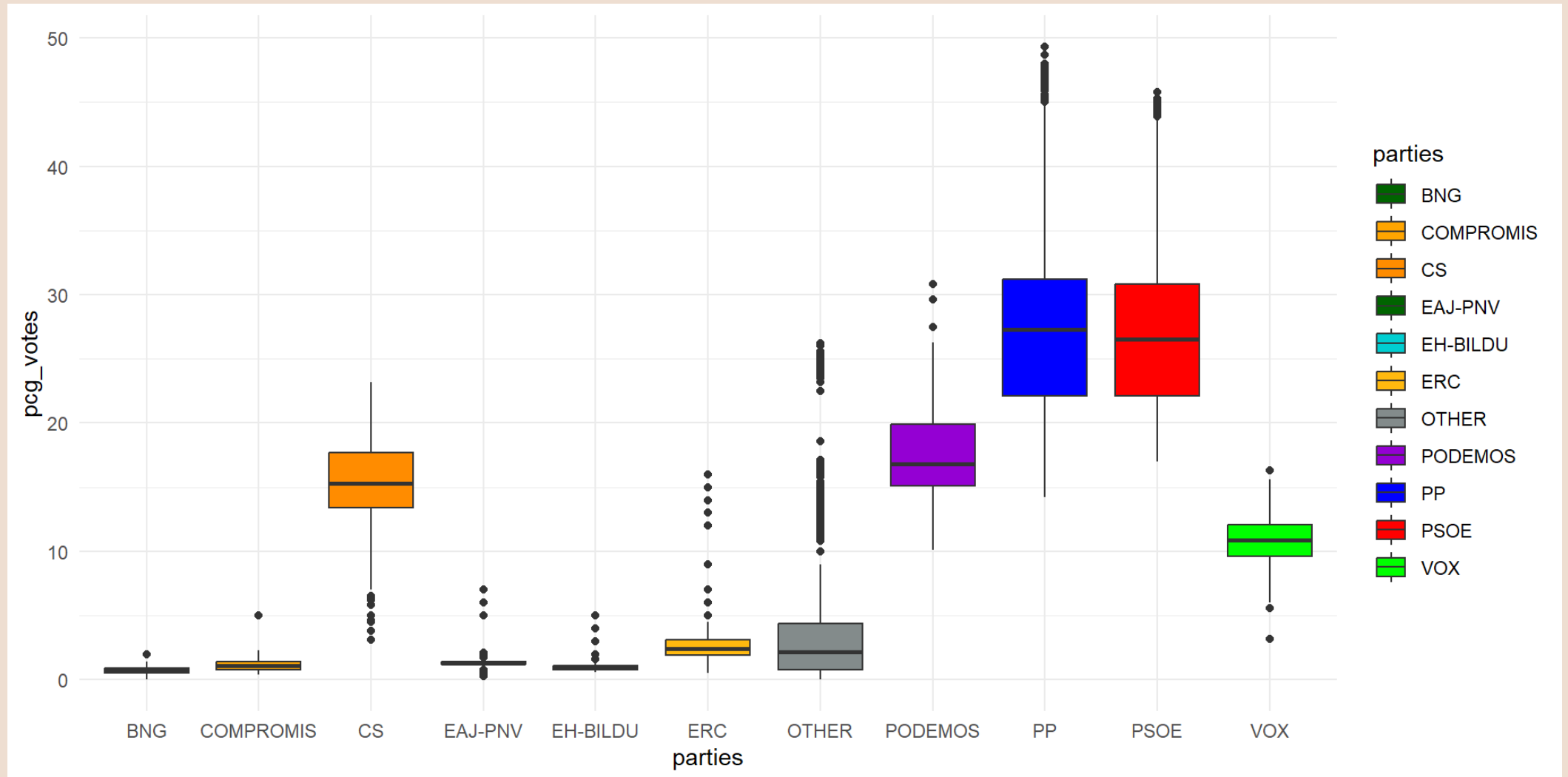
```
1 low_turnout <- surveys_pivot |> filter(low_turnout < 1 & votes_percent > 8.939)
2
3 low_turnout |>
4   distinct(parties)

# A tibble: 7 × 1
#   parties
#   <chr>
1 PSOE
2 PP
3 CS
4 OTHER
5 PODEMOS
6 VOX
7 ERC
```

These are the parties who benefited from low turnout, because they have more votes than expected.

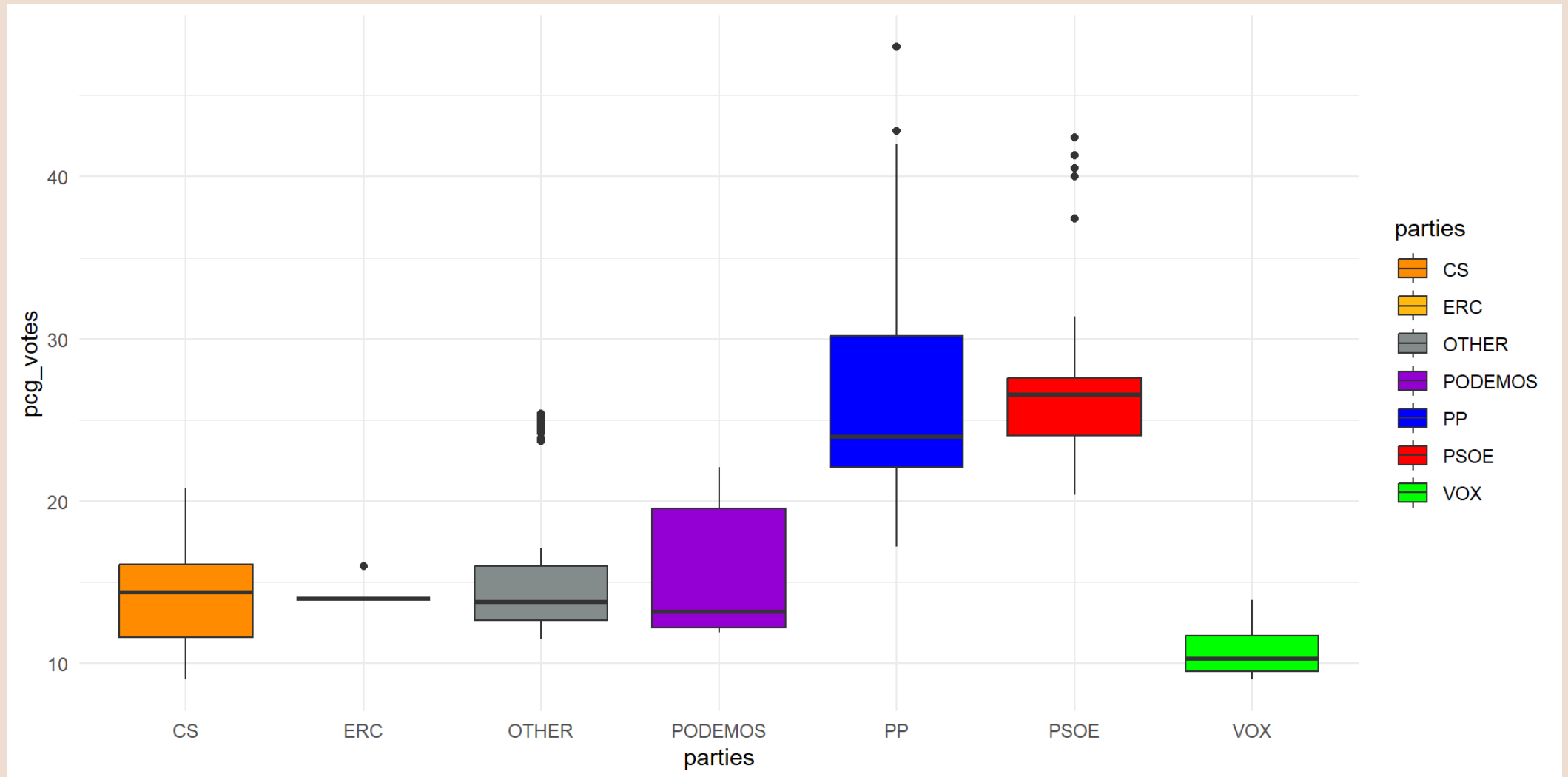
2.4 Who benefits from low turnout?

To visualize the results more clearly, boxplots have been created.



2.4 Who benefits from low turnout?

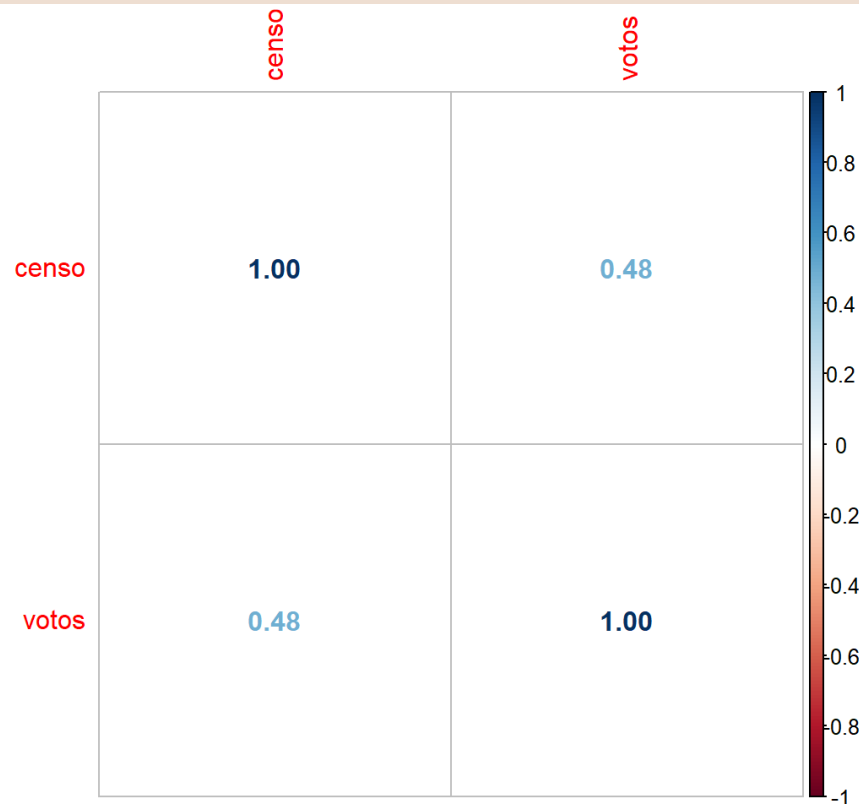
To visualize the results more clearly, boxplots have been created.



2.5 How to analyze the relationship between census and vote?

2.5 How to analyze the relationship between census and vote?

In order to analyze the relationship between the census and votes, a correlation plot has been generated to observe how these variables are correlated with each other. As we can see, they exhibit a moderate positive correlation of 0.48.



2.5 How to analyze the relationship between census and vote?

Additionally, a linear regression model has been created to determine if these variables are significant and may have an effect on each other.

```
Call:
lm(formula = censo ~ votos, data = election_pivot)

Residuals:
    Min       1Q   Median       3Q      Max
-2329601   -5130    -4688   -2521   2374697

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.188e+03  6.686e+01   77.6    <2e-16 ***
votos        4.988e+00  1.453e-02   343.2    <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 41970 on 396733 degrees of freedom
Multiple R-squared:  0.2289,    Adjusted R-squared:  0.2289
F-statistic: 1.178e+05 on 1 and 396733 DF,  p-value: < 2.2e-16
```


2.5 Is it true that certain parties win in rural areas?

As rural provinces, the following have been selected, as they are provinces in the depopulated Spain.

- 01-Álava 02-Albacete
- 05-Ávila 12-Castellón
- 13-Ciudad Real 16-Cuenca
- 19-Guadalajara 21-Huelva
- 22-Huesca 23-Jaén
- 25-Lleida, 26-La Rioja
- 27-Lugo 31-Navarra
- 32-Ourense 39-Cantabria,
- 40-Segovia 42-Soria
- 44-Teruel 49-Zamora

2.5 Is it true that certain parties win in rural areas?

```
1 depopulated_spain <- election_pivot |> filter(codigo_provincia == "01"|codigo_provincia == "02"|codigo_provincia == "03")
2
3 win_depop_spain <- depopulated_spain |> select(c(codigo_provincia, date_elec_ym, siglas, votos))
4 win_depop_spain
```

```
# A tibble: 140,144 × 4
```

| | codigo_provincia <fct> | date_elec_ym <date> | siglas <chr> | votos <dbl> |
|----|---------------------------|------------------------|-----------------|----------------|
| 1 | 01 | 2008-03-01 | PODEMOS | 9 |
| 2 | 01 | 2008-03-01 | OTHERS | 27 |
| 3 | 01 | 2008-03-01 | PSOE | 1 |
| 4 | 01 | 2008-03-01 | OTHERS | 1 |
| 5 | 01 | 2008-03-01 | OTHERS | 2 |
| 6 | 01 | 2008-03-01 | PP | 238 |
| 7 | 01 | 2008-03-01 | PODEMOS | 61 |
| 8 | 01 | 2008-03-01 | OTHERS | 85 |
| 9 | 01 | 2008-03-01 | OTHERS | 4 |
| 10 | 01 | 2008-03-01 | OTHERS | 17 |

```
# i 140,134 more rows
```

2.5 Is it true that certain parties win in rural areas?

Therefore, a new data frame named “maxvotes” has been created to consider the province code, election date, political party, and the number of votes each party received.

```
1 maxvotes <- win_depop_spain |>
2   filter(date_elec_ym %in% c("2008-03-01", "2011-11-01", "2015-12-01", "2016-06-01", "2019-04-01", "2019-11-01"))
3   group_by(date_elec_ym, codigo_provincia) |>
4   slice_max(votos, n=1)
5 maxvotes
```

```
# A tibble: 120 × 4
```

```
# Groups:   date_elec_ym, codigo_provincia [120]
```

| | codigo_provincia | date_elec_ym | siglas | votos |
|----|------------------|--------------|--------|-------|
| | <fct> | <date> | <chr> | <dbl> |
| 1 | 01 | 2008-03-01 | PSOE | 56349 |
| 2 | 02 | 2008-03-01 | PP | 49909 |
| 3 | 05 | 2008-03-01 | PP | 20468 |
| 4 | 12 | 2008-03-01 | PP | 42498 |
| 5 | 13 | 2008-03-01 | PP | 23826 |
| 6 | 16 | 2008-03-01 | PP | 15943 |
| 7 | 19 | 2008-03-01 | PP | 23910 |
| 8 | 21 | 2008-03-01 | PSOE | 37930 |
| 9 | 22 | 2008-03-01 | PSOE | 13227 |
| 10 | 23 | 2008-03-01 | PSOE | 33232 |

```
# i 110 more rows
```

2.5 Is it true that certain parties win in rural areas?

We can observe how in most provinces, the parties that won in the elections over the years have not always been the same. However, in Ávila, the PP wins over time.

```
1 alava_parties <- maxvotes |> filter(codigo_provincia == "01")
2 albacete_parties <- maxvotes |> filter(codigo_provincia == "02")
3 avila_parties <- maxvotes |> filter(codigo_provincia == "05")
4 avila_parties
```

```
# A tibble: 6 × 4
# Groups:   date_elec_ym, codigo_provincia [6]
  codigo_provincia date_elec_ym siglas votos
  <fct>            <date>      <chr>  <dbl>
1 05               2008-03-01    PP    20468
2 05               2011-11-01    PP    19391
3 05               2015-12-01    PP    14125
4 05               2016-06-01    PP    15553
5 05               2019-04-01    PP     9084
6 05               2019-11-01    PP     9260
```

2.5 Is it true that certain parties win in rural areas?

In Ciudad Real, the PP predominantly wins.

```
# A tibble: 6 × 4
# Groups:   date_elec_ym, codigo_provincia [6]
  codigo_provincia date_elec_ym siglas votos
  <fct>           <date>      <chr>  <dbl>
1 13              2008-03-01    PP    23826
2 13              2011-11-01    PP    24946
3 13              2015-12-01    PP    18080
4 13              2016-06-01    PP    19605
5 13              2019-04-01  PSOE    12838
6 13              2019-11-01    PP    12852
```

2.5 Is it true that certain parties win in rural areas?

In Cantabria, the PP predominantly wins.

```
# A tibble: 6 × 4
# Groups:   date_elec_ym, codigo_provincia [6]
  codigo_provincia date_elec_ym siglas votos
  <fct>           <date>      <chr>  <dbl>
1 39              2008-03-01    PP    58821
2 39              2011-11-01    PP    56866
3 39              2015-12-01    PP    39839
4 39              2016-06-01    PP    43755
5 39              2019-04-01  PSOE    25466
6 39              2019-11-01    PP    28019
```

2.5 Is it true that certain parties win in rural areas?

Lastly, in Segovia, the PP predominantly wins.

```
# A tibble: 6 × 4
# Groups:   date_elec_ym, codigo_provincia [6]
  codigo_provincia date_elec_ym siglas votos
  <fct>           <date>      <chr>  <dbl>
1 40              2008-03-01    PP    16555
2 40              2011-11-01    PP    15867
3 40              2015-12-01    PP    10669
4 40              2016-06-01    PP    12149
5 40              2019-04-01   PSOE     9286
6 40              2019-11-01    PP     9213
```

2.6 How to calibrate the error of the polls (remember that the polls are voting intentions at national level)?

2.6 How to calibrate the error of the polls?

We first need to convert the votes obtained by the parties in each election into percentages.

- Total votes by election

```
1 election_pivot_total <- election_pivot |>
2   group_by(date_elec_ym) |>
3   summarise(total_votes_election = sum(votos)) |>
4   print()
```

```
# A tibble: 6 × 2
  date_elec_ym total_votes_election
  <date>         <dbl>
1 2008-03-01      25069038
2 2011-11-01      23918052
3 2015-12-01      24895828
4 2016-06-01      23754401
5 2019-04-01      25877751
6 2019-11-01      23862002
```

2.6 How to calibrate the error of the polls (remember that the polls are voting intentions at national level)?

- Total votes obtained by each party in each election

```
1 election_pivot_votes <- election_pivot |>
2   group_by(date_elec_ym, siglas) |>
3   summarise(total_votes_party = sum(votos))
```

- Percentage

```
1 votes_election <-
2   inner_join(election_pivot_total , election_pivot_votes, by="date_elec_ym") |>
3   mutate(result_election = 100*(total_votes_party / total_votes_election))
```

2.6 How to calibrate the error of the polls (remember that the polls are voting intentions at national level)?

Since there are several polls for the same election, conducted by different media, we can compute the mean to obtain the average estimated result for each party and election.

```
1 votes_survey <- surveys_pivot |>
2   group_by(date_elec_ym, party) |>
3   summarize(result_poll = mean(votes_percent))
```

Once we have percentage in both election and poll data, we merge the results. Then, the difference between the poll estimation and the real results is computed in order to get the polling error.

```
1 merged_votes <-
2   inner_join(votes_election, votes_survey,
3             by = c("date_elec_ym" = "date_elec_ym",
4                   "siglas" = "party")) |>
5   mutate(poll_error = result_election - result_poll)
```

2.6 How to calibrate the error of the polls (remember that the polls are voting intentions at national level)?

```
1 print(merged_votes |>
2       select(date_elec_ym, siglas, result_election, result_poll, poll_error))
```

```
# A tibble: 39 × 5
  date_elec_ym siglas  result_election result_poll poll_error
  <date>      <chr>      <dbl>         <dbl>    <dbl>
1 2008-03-01   BNG          0.835         1.08     -0.248
2 2008-03-01  EAJ-PNV       1.21         1.79     -0.581
3 2008-03-01   ERC          1.18         2.33     -1.15
4 2008-03-01   PP          40.6         38.4      2.18
5 2008-03-01  PSOE         44.2         42.2      1.99
6 2011-11-01   BNG          0.766         1.08     -0.316
7 2011-11-01  EAJ-PNV       1.35         1.22      0.131
8 2011-11-01   ERC          1.07         0.962     0.110
9 2011-11-01   PP          45.3         43.3      2.05
10 2011-11-01  PSOE         29.2         36.0     -6.82
# i 29 more rows
```

2.7 In which election were the polls most wrong?

2.7 In which election were the polls most wrong?

To analyze which election had the most inaccurate polls, we first compute the average error for each election and search for the date with the highest value.

```
# A tibble: 6 × 2
  date_elec_ym avg_poll_error
  <date>         <dbl>
1 2008-03-01      1.23
2 2011-11-01      1.89
3 2015-12-01      1.23
4 2016-06-01      1.04
5 2019-04-01      0.78
6 2019-11-01      1.41
```

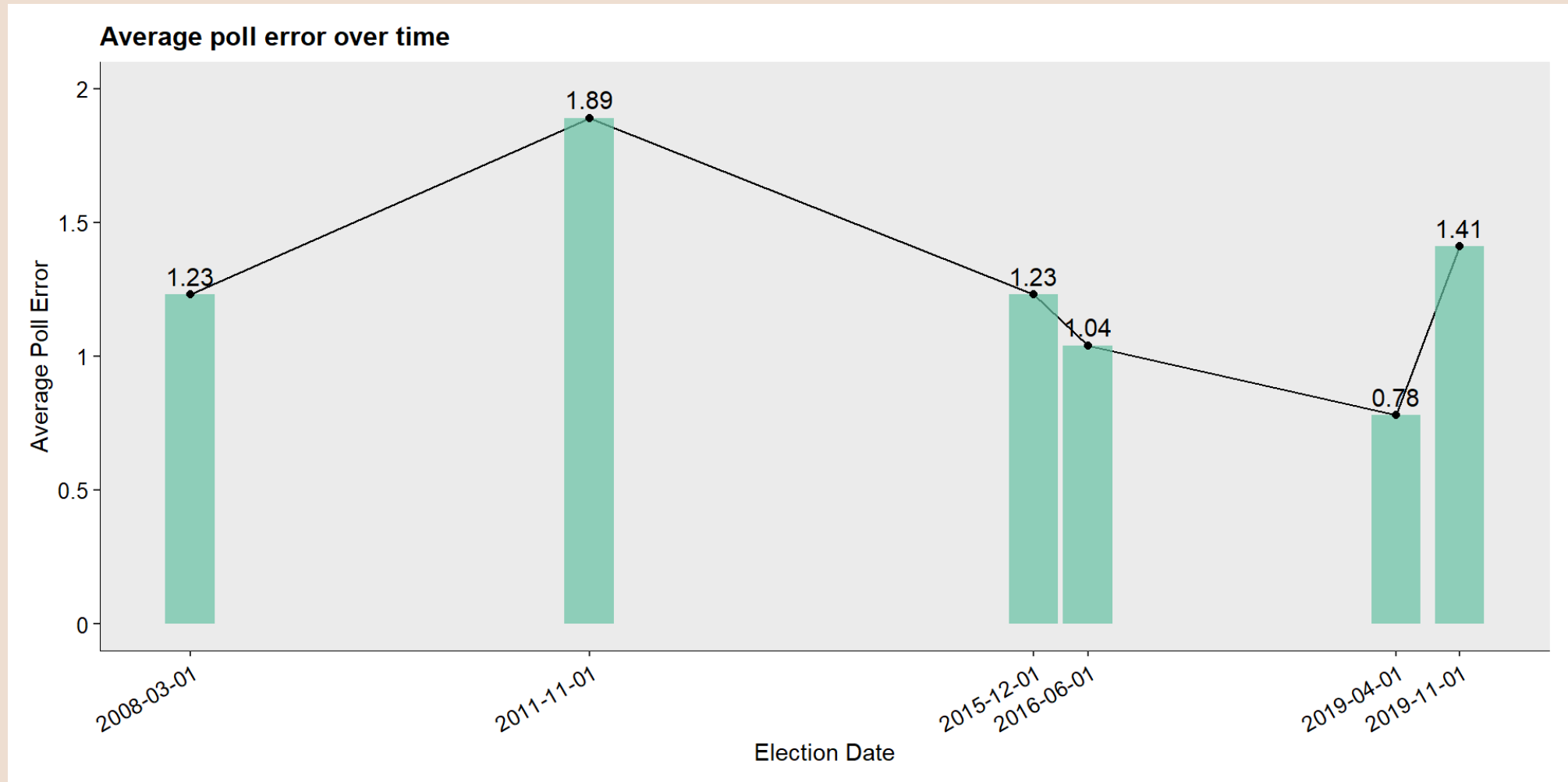
```
1 max_avg_error <- avg_error |>
2   slice_max(order_by = abs(avg_poll_error), n = 1) |>
3   print()
```

```
# A tibble: 1 × 2
  date_elec_ym avg_poll_error
  <date>         <dbl>
1 2011-11-01      1.89
```

The polls conducted on November 1, 2011, had the highest average error (1.89). Therefore, we can assume that these polls experienced the most significant inaccuracies, suggesting potential challenges or issues in their predictions.

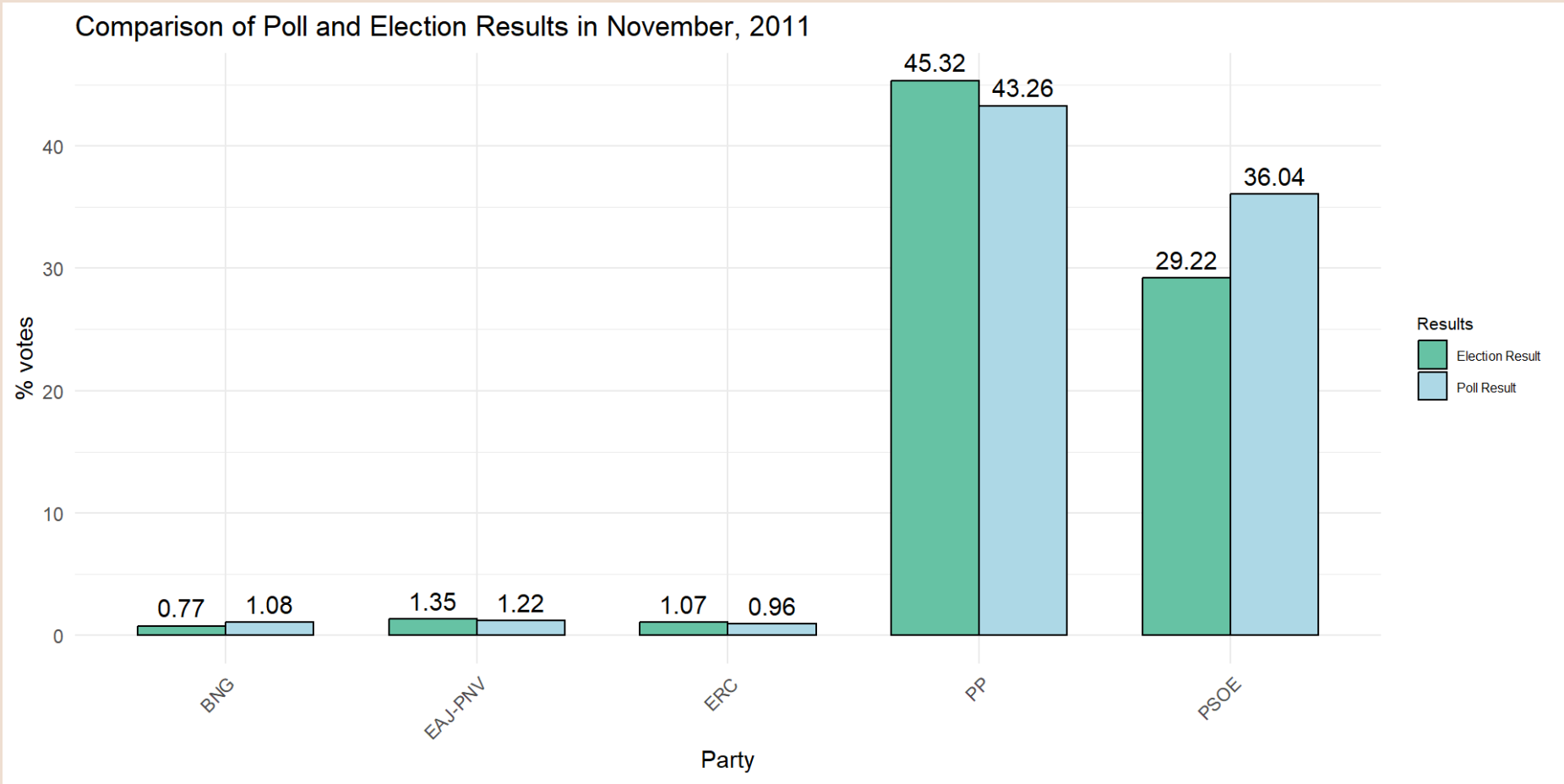
2.7 In which election were the polls most wrong?

The evolution of the average error can be represented for every election to observe the deviation in poll predictions over time.



2.7 In which election were the polls most wrong?

We further analyse the results obtained in these polls. By doing so, we can observe the deviation between the estimated results by the media and the real results obtained by the parties.



2.8 How were the polls wrong in national parties (PSOE, PP, VOX, CS, MP, UP - IU)?

2.8 How were the polls wrong in national parties (PSOE, PP, VOX, CS, MP, UP - IU)?

PSOE in 2011 was very over estimated by the polls

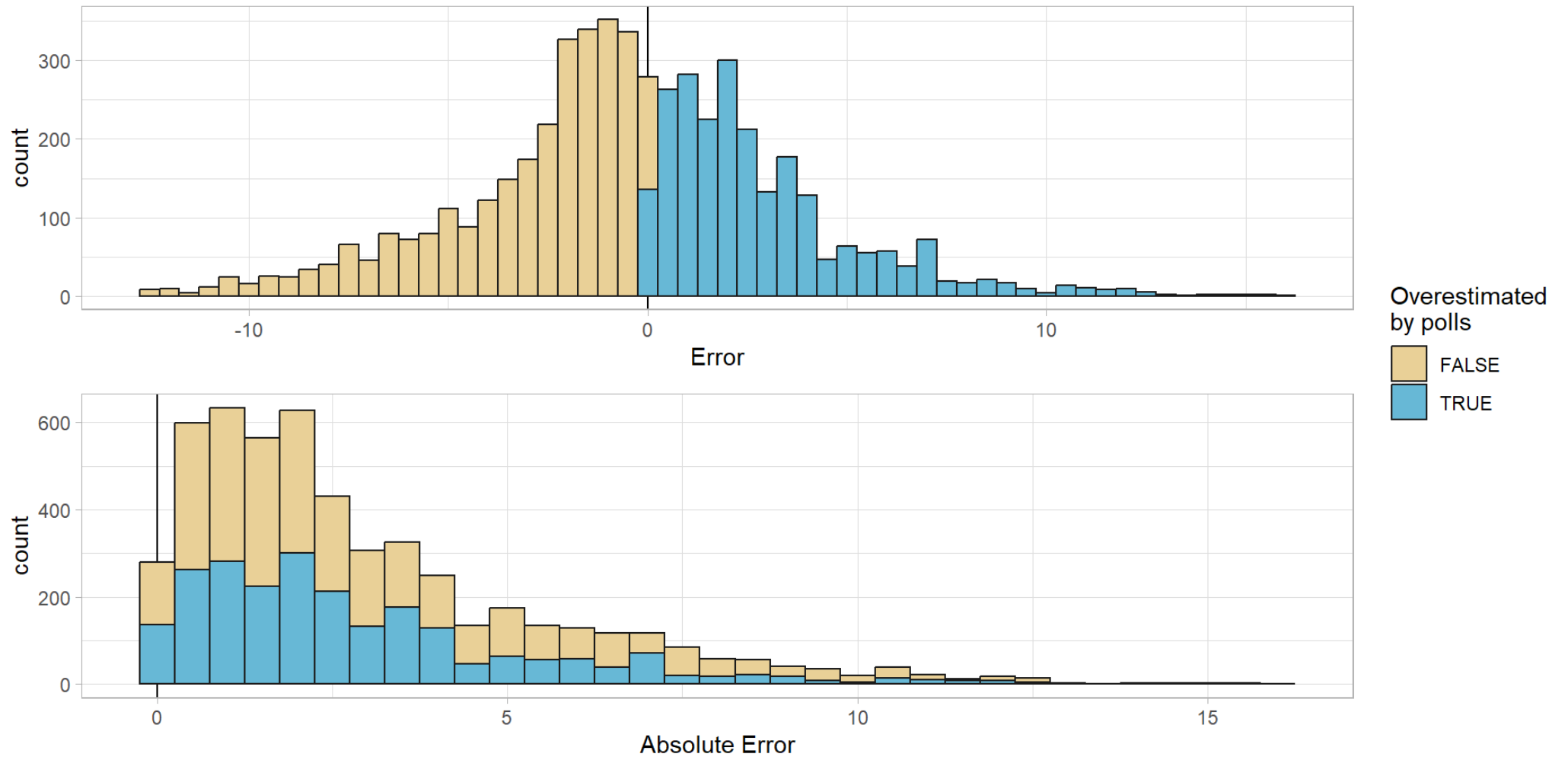
```
1 national_parties <- c("PSOE", "PP", "VOX", "CS", "M País", "PODEMOS")
2
3 merged_votes <- inner_join(votes_election, surveys_pivot,
4                             by = c("date_elec_ym" = "date_elec_ym",
5                                     "siglas" = "party"),
6                             multiple = "all") |>
7   mutate(error = votes_percent - result_election,
8           abs_error = abs(error),
9           overestimated = if_else(error > 0, T, F)
10          ) |>
11   filter(siglas %in% national_parties) |>
12   select(-c(low_turnout))
```

2.8 How were the polls wrong in national parties (PSOE, PP, VOX, CS, MP, UP - IU)?

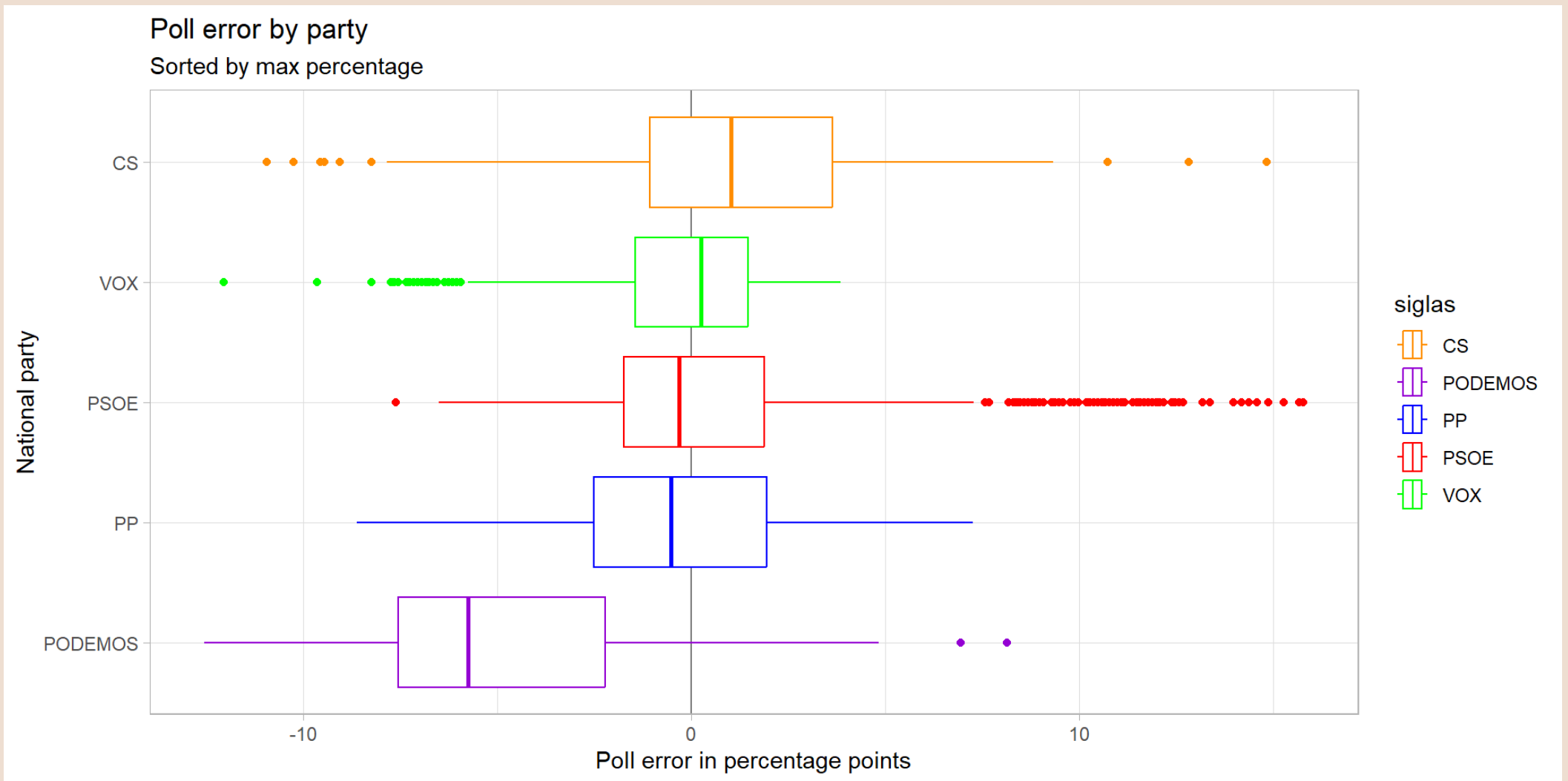
```
# A tibble: 10 × 6
```

| | date_elec | siglas | result_election | pollster | votes_percent | error |
|----|------------|--------|-----------------|---------------|---------------|-------|
| | <date> | <chr> | <dbl> | <chr> | <dbl> | <dbl> |
| 1 | 2011-11-20 | PSOE | 29.2 | GESOP | 45 | 15.8 |
| 2 | 2011-11-20 | PSOE | 29.2 | CELESTE-TEL | 44.9 | 15.7 |
| 3 | 2011-11-20 | PSOE | 29.2 | CELESTE-TEL | 44.5 | 15.3 |
| 4 | 2011-11-20 | PSOE | 29.2 | ASEP | 44.1 | 14.9 |
| 5 | 2019-11-10 | CS | 6.86 | SIMPLE LÓGICA | 21.7 | 14.8 |
| 6 | 2011-11-20 | PSOE | 29.2 | ASEP | 43.8 | 14.6 |
| 7 | 2011-11-20 | PSOE | 29.2 | CIS | 43.6 | 14.4 |
| 8 | 2011-11-20 | PSOE | 29.2 | CELESTE-TEL | 43.4 | 14.2 |
| 9 | 2011-11-20 | PSOE | 29.2 | SIMPLE LÓGICA | 43.2 | 14.0 |
| 10 | 2011-11-20 | PSOE | 29.2 | SIGMA DOS | 42.6 | 13.4 |

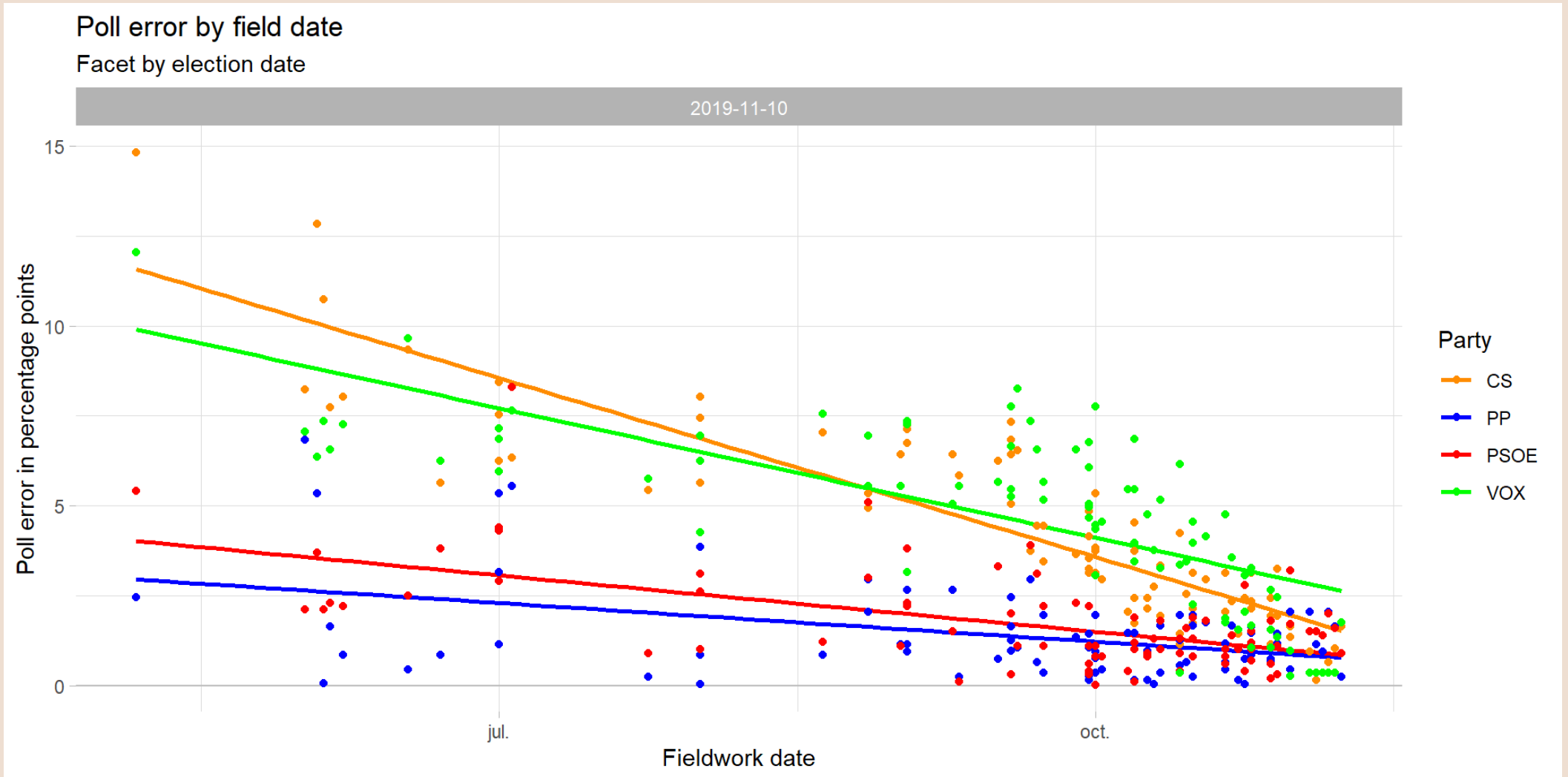
Poll error distribution



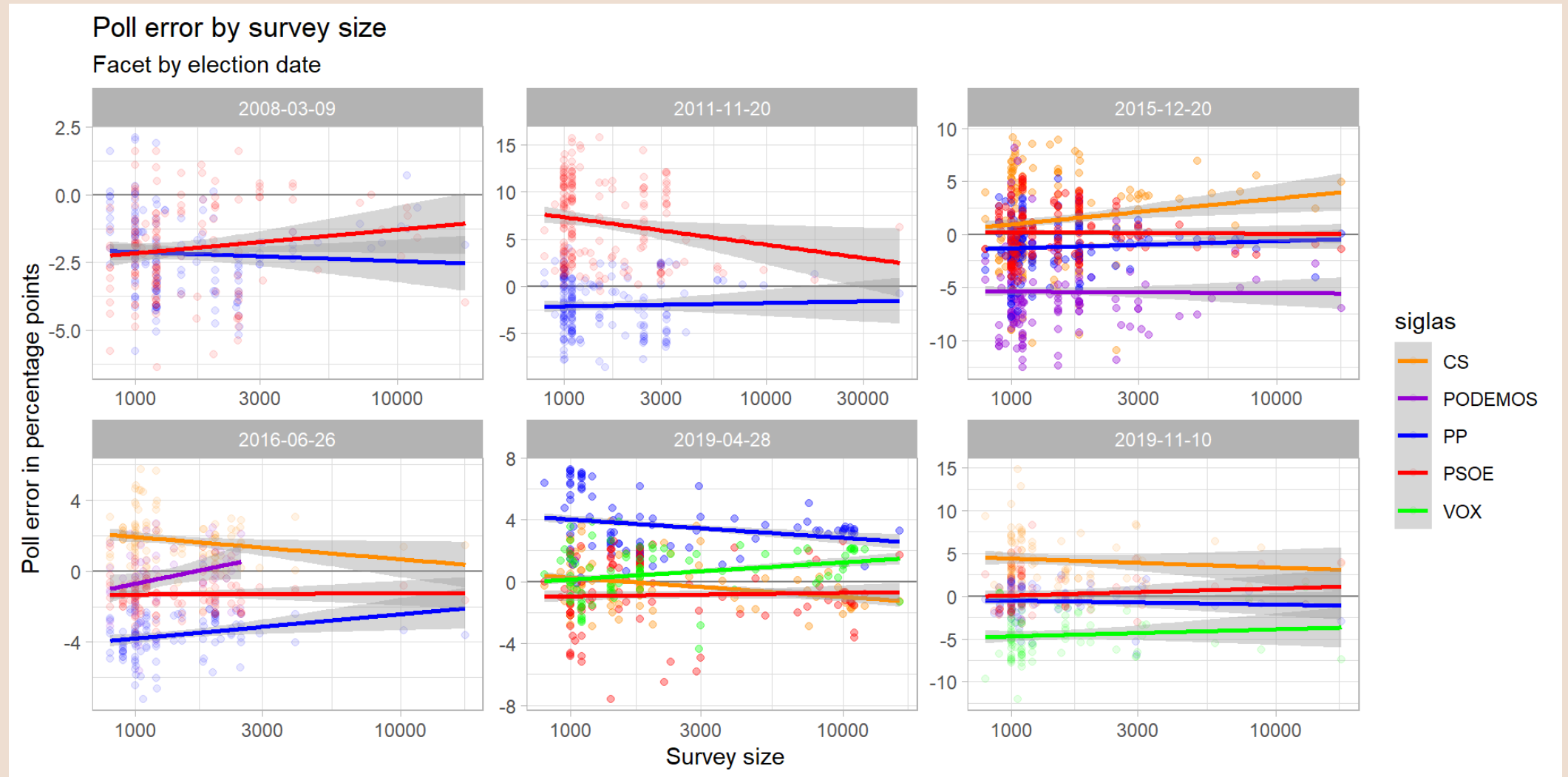
Most polls underestimated PODEMOS



Fieldwork date affected accuracy of Nov 2019 election



Did size have an affect on poll error?



2.9 Which polling houses got it right the most and which ones deviated the most from the results?

Each polling house surveyed a different number of parties

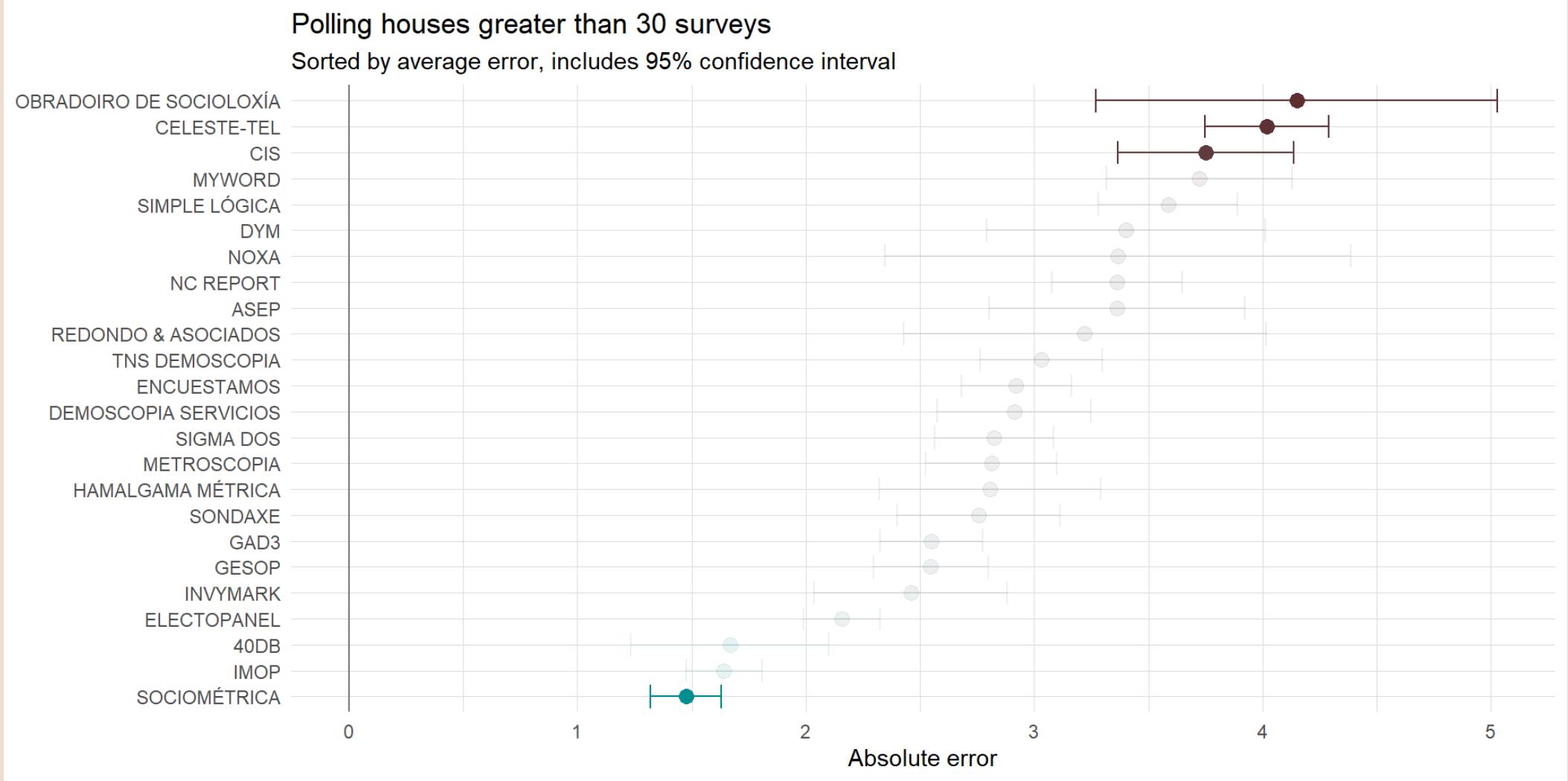
| | CS | PODEMOS | PP | PSOE | VOX |
|-----------------------------------|-----|---------|-----|------|-----|
| 40DB | 10 | 0 | 10 | 10 | 10 |
| A+M | 5 | 4 | 5 | 5 | 0 |
| ADVICE STRATEGIC | 3 | 1 | 3 | 3 | 0 |
| APPEND | 0 | 0 | 1 | 1 | 0 |
| ASEP | 0 | 0 | 39 | 39 | 0 |
| CELESTE-TEL | 117 | 68 | 164 | 164 | 43 |
| CEMOP | 0 | 0 | 1 | 1 | 0 |
| CIS | 39 | 22 | 70 | 70 | 16 |
| DEMOMÉTRICA | 0 | 0 | 6 | 6 | 0 |
| DEMOSCOPIA SERVICIOS | 45 | 32 | 45 | 45 | 13 |
| DYM | 16 | 10 | 28 | 28 | 6 |
| ELECTOPANEL | 80 | 0 | 80 | 80 | 80 |
| ENCUESTAMOS | 56 | 53 | 56 | 56 | 0 |
| ESTUDIO DE SOCIOLOGÍA CONSULTORES | 3 | 3 | 3 | 3 | 0 |
| GAD3 | 107 | 52 | 109 | 109 | 50 |
| GALLUP | 0 | 0 | 7 | 7 | 0 |
| GESOP | 51 | 21 | 75 | 75 | 16 |
| GIPEYOP | 7 | 5 | 7 | 7 | 0 |
| HAMATICAMA MÉTRICA | 12 | 4 | 12 | 12 | 0 |

Obradoiro de Socioloxía have the worst average error

```
1 subgroup <- merged_votes |>
2   group_by(pollster) |>
3   summarise(
4     avg_error = mean(abs_error),
5     var_error = var(abs_error),
6     sd_error = sd(abs_error),
7     total_surveys = n(),
8     num_parties = length(unique(siglas))
9   )
10
11 subgroup |> select(-var_error) |>
12   slice_max(avg_error, n = 5)
```

```
# A tibble: 5 × 5
  pollster                avg_error sd_error total_surveys num_parties
  <chr>                <dbl>   <dbl>         <int>         <int>
1 OBRADOIRO DE SOCIOLOXÍA  4.15    3.64            66             2
2 CELESTE-TEL            4.02    3.25           556             5
3 INTERCAMPO            4.01    3.36            12             2
4 A+M                   3.84    3.11            19             4
5 CIS                   3.75    2.89           217             5
```

2.9 Which polling houses got it right the most and which ones deviated the most from the results?



Which polling houses had the best estimates?

```
1 subgroup |> select(-var_error) |>
2   slice_min(avg_error, n = 7)
```

A tibble: 7 × 5

| | pollster <chr> | avg_error <dbl> | sd_error <dbl> | total_surveys <int> | num_parties <int> |
|---|-----------------------------------|--------------------|-------------------|------------------------|----------------------|
| 1 | TÁBULA V | 1.10 | 1.39 | 2 | 2 |
| 2 | SW DEMOSCOPIA | 1.11 | 0.784 | 4 | 4 |
| 3 | SYM CONSULTING | 1.15 | 1.72 | 3 | 3 |
| 4 | DEMOMÉTRICA | 1.21 | 1.11 | 12 | 2 |
| 5 | OPINA | 1.38 | 1.14 | 2 | 2 |
| 6 | ESTUDIO DE SOCIOLOGÍA CONSULTORES | 1.45 | 1.36 | 12 | 4 |
| 7 | SOCIOMÉTRICA | 1.47 | 1.24 | 248 | 4 |

