

POLI3148. Data Science in Politics and Public Administration

Problem Set 1+2 (15% + 15%)

Due: 2023-12-3 23:59

General Introduction

In this Problem Set, you will apply data science skills to wrangle and visualize the replication data of the following research article:

Cantú, F. (2019). The fingerprints of fraud: Evidence from Mexico's 1988 presidential election. *American Political Science Review*, 113(3), 710-726.

Requirements and Reminders

- You are required to use **RMarkdown** to compile your answer to this Problem Set.
- Use the attached compressed repo `POLI3148_PS_code_data.zip` to obtain starter code and data.
- Two submissions are required (via Moodle)
 - A `.pdf` file rendered by **Rmarkdown** that contains all your answer.
 - A compressed (in `.zip` format) R project repo. The expectation is that the instructor can unzip, open the project file, knitr your `.Rmd` file, and obtain the exact same output as the submitted `.pdf` document.
- The Problem Set is worth 30 points in total, allocated across 7 tasks. The point distribution across tasks is specified in the title line of each task. Within each task, the points are evenly distributed across sub-tasks. Bonus points (+5% max.) will be awarded to recognize exceptional performance.
- Grading rubrics: Overall, your answer will be evaluated based on its quality in three dimensions
 - Correctness and beauty of your outputs
 - Style of your code
 - Insightfulness of your interpretation or discussion
- Unless otherwise specified, you are required to use functions from the **tidyverse** package to complete this assignments.
- For some tasks, there may be multiple ways to achieve the same desired outcomes. You are encouraged to explore multiple methods. If you perform a task using multiple methods, do show it in your submission. You may earn bonus points for it.
- You are encouraged to use Generative AI such as ChatGPT to assist with your work. However, you will need to acknowledge it properly and validate AI's outputs. You may attach selected chat history with the AI you use and describe how it helps you get the work done. Extra credit may be rewarded to recognize creative use of Generative AI.
- This Problem Set is an individual assignment. You are expected to complete it independently. Clarification questions are welcome. Discussions on concepts and techniques related to the Problem Set among peers is encouraged. However, without the instructor's consent, sharing (sending and requesting) code and text that complete the entirety of a task is prohibited. You are strongly encouraged to use *CampusWire* for clarification questions and discussions.

Background

In 1998, Mexico had a close presidential election. Irregularities were detected around the country during the voting process. For example, when 2% of the vote tallies had been counted, the preliminary results showed the PRI's imminent defeat in Mexico City metropolitan area and a very narrow vote margin between PRI and FDN. A few minutes later, the screens at the Ministry of Interior went blank, an event that electoral authorities justified as a technical problem caused by an overload on telephone lines. The vote count was therefore suspended for three days, despite the fact that opposition representatives found a computer in the basement that continued to receive electoral results. Three days later, the vote count resumed, and soon the official announced PRI's winning with 50.4% of the vote.

What happened on that night and the following days? Were there electoral fraud during the election? A political scientist, Francisco Cantú, unearths a promising dataset that could provide some clues. At the National Archive in Mexico City, Cantú discovered about 53,000 vote tally sheets. Using machine learning methods, he detected that a significant number of tally sheets were *altered*! In addition, he found evidence that the altered tally sheets were biased in favor of the incumbent party. In this Problem Set, you will use Cantú's replication dossier to replicate and extend his data work.

Please read Cantú (2019) for the full story. And see Figure 1 for a few examples of altered (fraudulent) tallies.

VOTACION RECIBIDA EN LA URNA (con número)	VOTOS ENCONTRADOS EN OTRAS URNAS (con número)	(con número)
131		131
97		7
138		138
138		138
138		138
138		138

VOTACION RECIBIDA EN LA URNA (con número)	VOTOS ENCONTRADOS EN OTRAS URNAS (con número)	(con número)
29		
120		
131		
1		
10		
37		
1		
22		
2		
273		
14		
287		

VOTACION RECIBIDA EN LA URNA (con número)	VOTOS ENCONTRADOS EN OTRAS URNAS (con número)	(con número)
12		
1399		
20		
1		
2		
3		
1437		
1		
1437		

VOTACION RECIBIDA EN LA URNA (con número)	VOTOS ENCONTRADOS EN OTRAS URNAS (con número)	(con número)
359		359
22		22
381		381
381		381

Figure 1: Examples of altered tally sheets (reproducing Figure 1 of Cantú 2018)

Task 0. Loading required packages (3pt)

For Better organization, it is a good habit to load all required packages up front at the start of your document. Please load the all packages you use throughout the whole Problem Set here.

```
# YOUR CODE HERE
```

Task 1. Clean machine classification results (3pt)

Cantú applies machine learning models to 55,334 images of tally sheets to detect signs of fraud (i.e., alteration). The machine learning model returns results recorded in a table. The information in this table is messy and requires data wrangling before we can use them.

Task 1.1. Load classified images of tally sheets

The path of the classified images of tally sheets is `data/classification.txt`. Your first task is loading these data onto R using a `tidyverse` function. Name it `d_tally`.

Note:

- Although the file extension of this dataset is `.txt`, you are recommended to use the `tidyverse` function we use for `.csv` files to read it.
- Unlike the data files we have read in class, this table has *no column names*. Look up the documentation and find a way to handle it.
- There will be three columns in this dataset, name them `name_image`, `label`, and `probability`.

Print your table to show your output.

```
# YOUR CODE HERE
```

The printed output should look like the following:

```
## # A tibble: 55,334 x 3
##   name_image                label probability
##   <chr>                  <chr> <chr>
## 1 Aguascalientes_I_2014-05-26 00.00.10.jpg [[0]] [[ 0.99919599]]
## 2 Aguascalientes_I_2014-05-26 00.00.17.jpg [[0]] [[ 0.95722806]]
## 3 Aguascalientes_I_2014-05-26 00.00.25.jpg [[0]] [[ 0.57690716]]
## 4 Aguascalientes_I_2014-05-26 00.00.31.jpg [[0]] [[ 0.96505082]]
## 5 Aguascalientes_I_2014-05-26 00.00.38.jpg [[0]] [[ 0.86975688]]
## 6 Aguascalientes_I_2014-05-26 00.00.45.jpg [[0]] [[ 0.78825063]]
## 7 Aguascalientes_I_2014-05-26 00.00.52.jpg [[0]] [[ 0.96493018]]
## 8 Aguascalientes_I_2014-05-26 00.00.59.jpg [[0]] [[ 0.68087846]]
## 9 Aguascalientes_I_2014-05-26 00.01.06.jpg [[0]] [[ 0.99999994]]
## 10 Aguascalientes_I_2014-05-26 00.01.15.jpg [[0]] [[ 0.64047635]]
## # i 55,324 more rows
```

Hint: Please make sure the number of rows, the number of columns, and the column names match with the above output.

Note 1. What are in this dataset?

Before you proceed, let me explain the meaning of the three variables.

- **name_image** contains the names of the tallies' image files (as you may infer from the .jpg file extensions. They contain information about the locations where each of the tally sheets are produced.
- **label** is a machine-predicted label indicating whether a tally is fraudulent or not. **label** = 1 means the machine learning model has detected signs of fraud in the tally sheet. **label** = 0 means the machine detects no sign of fraud in the tally sheet. In short, **label** = 1 means fraud; **label** = 0 means no fraud.
- **probability** indicates the machine's certainty about its predicted **label** (explained above). It ranges from 0 to 1, where higher values mean higher level of certainty.

Interpret **label** and **probability** carefully. Two examples can hopefully give you clues about their correct interpretation. In the first row, **label** = 0 and **probability** = 0.9991. That means the machine thinks this tally sheet is NOT FRAUDULENT with a probability of 0.9991. Then, the probability that this tally sheet is fraudulent is $1 - 0.9991 = 0.0009$. Take another example, in the 11th row, **label** = 1 and **probability** = 0.935. This means the machine thinks this tally sheet IS FRAUDULENT with a probability of 0.935. Then, the probability that it is NOT FRAUDULENT is $1 - 0.9354 = 0.0646$.

Task 1.2. Clean columns label and probability

As you have seen in the printed outputs, columns `label` and `probability` are read as `chr` variables when they are actually numbers. A close look at the data may tell you why — they are “wrapped” by some non-numeric characters. In this task, you will clean these two variables and make them valid numeric variables. You are required to use `tidyverse` operations to for this task. Show appropriate summary statistics of `label` and `probability` respectively after you have transformed them into numeric variables.

```
# YOUR CODE HERE
```

```
## # A tibble: 55,334 x 3
##   name_image          label probability
##   <chr>          <dbl>      <dbl>
## 1 Aguascalientes_I_2014-05-26 00.00.10.jpg      0      0.999
## 2 Aguascalientes_I_2014-05-26 00.00.17.jpg      0      0.957
## 3 Aguascalientes_I_2014-05-26 00.00.25.jpg      0      0.577
## 4 Aguascalientes_I_2014-05-26 00.00.31.jpg      0      0.965
## 5 Aguascalientes_I_2014-05-26 00.00.38.jpg      0      0.870
## 6 Aguascalientes_I_2014-05-26 00.00.45.jpg      0      0.788
## 7 Aguascalientes_I_2014-05-26 00.00.52.jpg      0      0.965
## 8 Aguascalientes_I_2014-05-26 00.00.59.jpg      0      0.681
## 9 Aguascalientes_I_2014-05-26 00.01.06.jpg      0      1.00
## 10 Aguascalientes_I_2014-05-26 00.01.15.jpg      0      0.640
## # i 55,324 more rows
```

Hints:

- You should be able to find the function you need in this *tidyverse* cheat sheet: <https://rstudio.github.io/cheatsheets/html/strings.html>
- `[` is a special character used in R's regular expression. When you refer to it in R, you need to add two `/` in front of it. For example
- You can do the above operation in multiple steps.

Task 1.3. Extract state and district information from name_image

As explained in the note, the column `name_image`, which has the names of tally sheets' images, contains information about locations where the tally sheets are produced. Specifically, the first two elements of these file names indicates the **states'** and districts' identifiers respectively, for example, `name_image = "Aguascalientes_I_2014-05-26 00.00.10.jpg"`. It means this tally sheet is produced in state **Aguascalientes**, district **I**. In this task, you are required to obtain this information. Specifically, create two columns named `state` and `district` as state and district identifiers respectively. You are required to use `tidyverse` functions to perform the task.

```
# YOUR CODE HERE
```

Your output should look like the following:

```
## # A tibble: 55,334 x 5
##   name_image                state    district label probability
##   <chr>                <chr>    <chr>    <dbl>      <dbl>
## 1 Aguascalientes_I_2014-05-26 00.00.10.jpg Aguascal~ I         0         0.999
## 2 Aguascalientes_I_2014-05-26 00.00.17.jpg Aguascal~ I         0         0.957
## 3 Aguascalientes_I_2014-05-26 00.00.25.jpg Aguascal~ I         0         0.577
## 4 Aguascalientes_I_2014-05-26 00.00.31.jpg Aguascal~ I         0         0.965
## 5 Aguascalientes_I_2014-05-26 00.00.38.jpg Aguascal~ I         0         0.870
## 6 Aguascalientes_I_2014-05-26 00.00.45.jpg Aguascal~ I         0         0.788
## 7 Aguascalientes_I_2014-05-26 00.00.52.jpg Aguascal~ I         0         0.965
## 8 Aguascalientes_I_2014-05-26 00.00.59.jpg Aguascal~ I         0         0.681
## 9 Aguascalientes_I_2014-05-26 00.01.06.jpg Aguascal~ I         0         1.00
## 10 Aguascalientes_I_2014-05-26 00.01.15.jpg Aguascal~ I         0         0.640
## # i 55,324 more rows
```

Hints:

- You may try the `separate` function.
- Alternatively, you may look up functions in the `stringr` module.

Task 1.4. Re-code a state's name

One of the states (in the newly created column `state`) is coded as “Estado de Mexico.” The researchers decide that it should instead re-coded as “**Edomex**.” Please use a `tidyverse` function to perform this task.

Hint: Look up functions `ifelse` and `case_match`.

```
# YOUR CODE HERE
```


Task 1.5. Create a *probability of fraud* indicator

As explained in Note 1, we need to interpret `label` and `probability` with caution, as the meaning of `probability` is conditional on the value of `label`. To avoid confusion in the analysis, your next task is to create a column named `fraud_proba` which indicates the probability that a tally sheet is fraudulent. After you have created the column, drop the `label` and `probability` columns.

Hint: Look up the `ifelse` function and the `case_when` function (but you just need either one of them).

```
# YOUR CODE HERE
```

Your output should look like the following:

```
## # A tibble: 55,334 x 4
##   name_image                state    district fraud_proba
##   <chr>                <chr>    <chr>      <dbl>
## 1 Aguascalientes_I_2014-05-26 00.00.10.jpg Aguascalientes I      0.000804
## 2 Aguascalientes_I_2014-05-26 00.00.17.jpg Aguascalientes I      0.0428
## 3 Aguascalientes_I_2014-05-26 00.00.25.jpg Aguascalientes I      0.423
## 4 Aguascalientes_I_2014-05-26 00.00.31.jpg Aguascalientes I      0.0349
## 5 Aguascalientes_I_2014-05-26 00.00.38.jpg Aguascalientes I      0.130
## 6 Aguascalientes_I_2014-05-26 00.00.45.jpg Aguascalientes I      0.212
## 7 Aguascalientes_I_2014-05-26 00.00.52.jpg Aguascalientes I      0.0351
## 8 Aguascalientes_I_2014-05-26 00.00.59.jpg Aguascalientes I      0.319
## 9 Aguascalientes_I_2014-05-26 00.01.06.jpg Aguascalientes I      0.0000000600
## 10 Aguascalientes_I_2014-05-26 00.01.15.jpg Aguascalientes I      0.360
## # i 55,324 more rows
```

Task 1.6. Create a binary *fraud* indicator

In this task, you will create a binary indicator called `fraud_bin` indicating whether a tally sheet is fraudulent. Following the researcher's rule, we consider a tally sheet fraudulent only when the machine thinks it is at least 2/3 likely to be fraudulent. That is, `fraud_bin` is set to `TRUE` when `fraud_proba` is greater to 2/3 and is `FALSE` otherwise.

```
# YOUR CODE HERE
```

Your output should look like the following:

```
## # A tibble: 55,334 x 5
##   name_image                state district fraud_proba fraud_bin
##   <chr>                <chr> <chr>      <dbl> <lgl>
## 1 Aguascalientes_I_2014-05-26 00.00.10.jpg Agua- I      8.04e-4 FALSE
## 2 Aguascalientes_I_2014-05-26 00.00.17.jpg Agua- I      4.28e-2 FALSE
## 3 Aguascalientes_I_2014-05-26 00.00.25.jpg Agua- I      4.23e-1 FALSE
## 4 Aguascalientes_I_2014-05-26 00.00.31.jpg Agua- I      3.49e-2 FALSE
## 5 Aguascalientes_I_2014-05-26 00.00.38.jpg Agua- I      1.30e-1 FALSE
## 6 Aguascalientes_I_2014-05-26 00.00.45.jpg Agua- I      2.12e-1 FALSE
## 7 Aguascalientes_I_2014-05-26 00.00.52.jpg Agua- I      3.51e-2 FALSE
## 8 Aguascalientes_I_2014-05-26 00.00.59.jpg Agua- I      3.19e-1 FALSE
## 9 Aguascalientes_I_2014-05-26 00.01.06.jpg Agua- I      6.00e-8 FALSE
## 10 Aguascalientes_I_2014-05-26 00.01.15.jpg Agua- I      3.60e-1 FALSE
## # i 55,324 more rows
```

Task 2. Visualize machine classification results (3pt)

In this section, you will visualize the `tally` dataset that you have cleaned in Task 1. Unless otherwise specified, you are required to use the `ggplot` packages to perform all the tasks.

Task 2.1. Visualize distribution of `fraud_proba`

How is the predicted probability of fraud (`fraud_proba`) distributed? Use two methods to visualize the distribution. Remember to add informative labels to the figure. Describe the plot with a few sentences.

```
# YOUR CODE HERE
```

Task 2.2. Visualize distribution of `fraud_bin`

How many tally sheets are fraudulent and how many are not? We may answer this question by visualizing the binary indicator of tally-level states of fraud. Use at least two methods to visualize the distribution of `fraud_bin`. Remember to add informative labels to the figure. Describe your plots with a few sentences.

```
# YOUR CODE HERE
```

The figure below serve as a reference. Feel free to try alternative approach(es) to make your visualization nicer and more informative.

Task 2.3. Summarize prevalence of fraud by state

Next, we will examine the between-state variation with regards to the prevalence of election fraud. In this task, you will create a new object that contains two state-level indicators regarding the prevalence of election fraud: The count of fraudulent tallies and the proportion of fraudulent tallies.

```
# YOUR CODE HERE
```

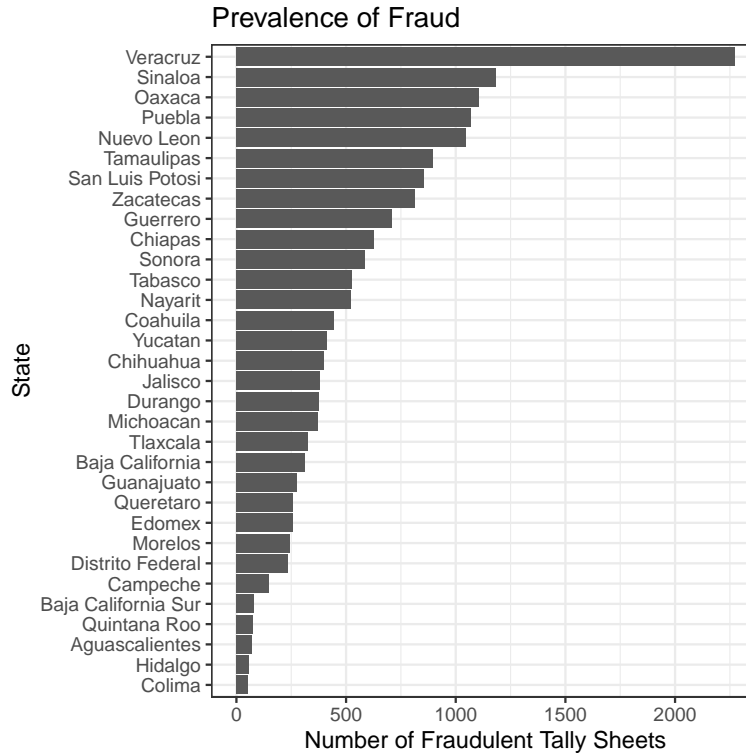
Your output should look like the following:

```
## # A tibble: 32 x 3
##   state          n_fraud prop_fraud
##   <chr>          <int>     <dbl>
## 1 Aguascalientes      71      17.6
## 2 Baja California    311      23.1
## 3 Baja California Sur   79      19.1
## 4 Campeche          146      38.6
## 5 Chiapas           629      45.6
## 6 Chihuahua          398      21.4
## 7 Coahuila           444      37.8
## 8 Colima              51      16.8
## 9 Distrito Federal   236       3.10
## 10 Durango            376      27.8
## # i 22 more rows
```

Task 2.4. Visualize frequencies of fraud by state

Using the new data frame created in Task 2.3, please visualize the *frequencies* of fraudulent tallies of every state. Describe the key takeaway from the visualization with a few sentences.

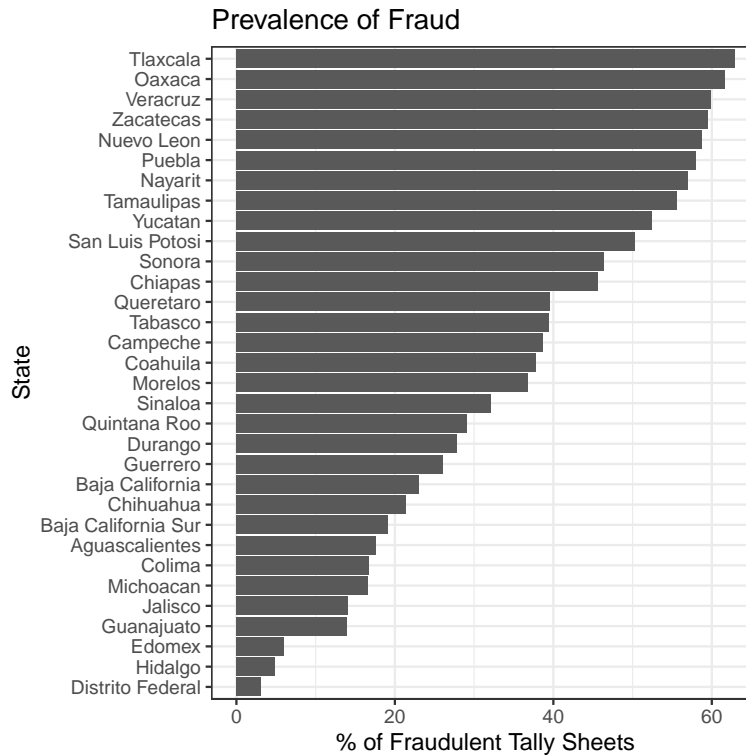
The figure below serve as a reference. Feel free to try alternative approach(es) to make your visualization nicer and more informative.



Task 2.5. Visualize proportions of fraud by state

Using the new data frame created in Task 2.3, please visualize the *proportion of* fraudulent tallies of every state. Describe the key takeaway from the visualization with a few sentences.

The figure below serve as a reference. Feel free to use alternative approach(es) to make your visualization more informative.



Task 2.6. Visualize both proportions & frequencies of fraud by state

Create data visualization to show BOTH the *proportions* and *frequencies* of fraudulent tally sheets by state in one figure. Include annotations to highlight states with the highest level of fraud. Add informative labels to the figure. Describe the takeaways from the figure with a few sentences.

```
# YOUR CODE HERE
```


Task 3. Clean vote return data (3pt)

Your next task is to clean a different dataset from the researchers' replication dossier. Its path is `data/Mexican_Election_Fraud/dataverse/VoteReturns.csv`. This dataset contains information about vote returns recorded in every tally sheet. This dataset is essential for the replication of Figure 4 in the research article.

Task 3.1. Load vote return data

Load the dataset onto your R environment. Name this dataset `d_return`. Show summary statistics of this dataset and describe the takeaways using a few sentences.

```
# YOUR CODE
```

Hint: the loaded object should look like the following:

```
## # A tibble: 53,499 x 91
##   foto seccion casilla dtto   dto municipio edo  entidad pagina  p1  p2
##   <chr> <chr>   <chr> <chr> <dbl> <chr>   <chr> <chr>   <dbl> <dbl> <dbl>
## 1 2014-- 83      83      I      1 AGUASCAL~ Agua~ AGS      127  108  333
## 2 2014-- 1       84     <NA>    1 AGUASCAL~ Agua~ AGUASC~  128  919  453
## 3 2014-- 85      85      1      1 AGUASCAL~ Agua~ AGUASC~  129  795  264
## 4 2014-- 45     45-A    1      1 AGUASCAL~ Agua~ AGUA    130  767  450
## 5 2014-- 86      86      1      1 AGUASCAL~ Agua~ AGUAS   131 1243  578
## 6 2014-- 87      87      1      1 <NA>      Agua~ 1      132  718  333
## 7 2014-- 1       87-A    7      1 AGUASCAL~ Agua~ AGUAS   133  710  299
## 8 2014-- 88      88      1      1 AGUAS      Agua~ AGUAS   134    0    0
## 9 2014-- 89      89      1      1 AGUASCAL~ Agua~ AGUAS   135  764    8
## 10 2014-- 89     89-A    7      1 AGUSCALI~ Agua~ 1      136  759  256
## # i 53,489 more rows
## # i 80 more variables: p3 <dbl>, p4 <dbl>, p5 <dbl>, pan <dbl>, pri <dbl>,
## #   pps <dbl>, psm <dbl>, pms <dbl>, pfcrrn <dbl>, prt <dbl>, parm <dbl>,
## #   noregis <dbl>, nombrenore <chr>, otros <dbl>, otroscan <chr>, pan2 <dbl>,
## #   pri2 <dbl>, pps2 <dbl>, psm2 <dbl>, pms2 <dbl>, pfcrrn2 <dbl>, prt2 <dbl>,
## #   parm2 <dbl>, noregis2 <dbl>, otro2 <dbl>, pan3 <dbl>, pri3 <dbl>,
## #   pps3 <dbl>, psm3 <dbl>, pms3 <dbl>, pfcrrn3 <dbl>, prt3 <dbl>, ...
```

Note 2. What are in this dataset?

This table contains a lot of different variables. The researcher offers no comprehensive documentation to tell us what every column means. For the sake of this problem set, you only need to know the meanings of the following columns:

- **foto** is an identifier of the images of tally sheets in this dataset. We will need it to merge this dataset with the **d_tally** data.
- **edo** contains the names of states.
- **dto** contains the names of districts (in Arabic numbers).
- **salinas**, **clouthier**, and **ibarra** contain the counts of votes (as recorded in the tally sheets) for presidential candidates Salinas (PRI), Cardenas (FDN), and Clouthier (PAN). In addition, the summation of all three makes the total number of **presidential votes**.
- **total** contains the total number of **legislative votes**.

Task 3.2. Recode names of states

A state whose name is `Chihuahua` is mislabelled as `Chihuhua`. A state whose name is currently `Edomex` needs to be recoded to `Estado de Mexico`. Please re-code the names of these two states accordingly.

```
# YOUR CODE
```

Task 3.3. Recode districts' identifiers

Compare how districts' identifiers are recorded differently in the tally (**d_tally**) from vote return (**d_return**) datasets. Specifically, in the **d_tally** dataset, **district** contains Roman numbers while in the **d_return** dataset, **dto** contains Arabic numbers. Recode districts' identifiers in the **d_return** dataset to match those in the **d_tally** dataset. To complete this task, first summarize the values of the two district identifier columns in the two datasets respectively to verify the above claim. Then do the requested conversion.

Task 3.4. Create a name_image identifier for the d_return dataset

In the `d_return` dataset, create a column named `name_image` as the first column. The column concatenate values in the three columns: `edo`, `dto`, and `foto` with an underscore `_` as separators.

```
# YOUR CODE HERE
```

Below the post-wrangling columns `name_image`, `edo`, `dto`, and `foto` are printed (with other columns omitted). You may use this output to infer the type of operations you are required to do.

```
## # A tibble: 53,499 x 4
##   name_image          edo          dto      foto
##   <chr>             <chr>      <roman> <chr>
## 1 Aguascalientes_I_2014-05-26 00.00.04.JPG Aguascalientes I      2014-05-26 0~
## 2 Aguascalientes_I_2014-05-26 00.00.10    Aguascalientes I      2014-05-26 0~
## 3 Aguascalientes_I_2014-05-26 00.00.17    Aguascalientes I      2014-05-26 0~
## 4 Aguascalientes_I_2014-05-26 00.00.25    Aguascalientes I      2014-05-26 0~
## 5 Aguascalientes_I_2014-05-26 00.00.31    Aguascalientes I      2014-05-26 0~
## 6 Aguascalientes_I_2014-05-26 00.00.38    Aguascalientes I      2014-05-26 0~
## 7 Aguascalientes_I_2014-05-26 00.00.45    Aguascalientes I      2014-05-26 0~
## 8 Aguascalientes_I_2014-05-26 00.00.52    Aguascalientes I      2014-05-26 0~
## 9 Aguascalientes_I_2014-05-26 00.00.59    Aguascalientes I      2014-05-26 0~
## 10 Aguascalientes_I_2014-05-26 00.01.06    Aguascalientes I      2014-05-26 0~
## # i 53,489 more rows
```

Task 3.5. Wrangle the name_image column in two datasets

As a final step before merging `d_return` and `d_tally`, you are required to perform the following data wrangling. For the `name_image` column in BOTH `d_return` and `d_tally`:

- Convert all characters to lower case.
- Remove ending substring `.jpg`.

```
# YOUR CODE HERE
```

Task 3.6 Join classification results and vote returns

After you have successfully completed all the previous steps, join `d_return` and `d_tally` by column `name_image`. This task contains two part. First, use appropriate `tidyverse` functions to answer the following questions:

- How many rows are in `d_return` but not in `d_tally`? Which states and districts are they from?
- How many rows are in `d_tally` but not in `d_return`? Which states and districts are they from?

YOUR CODE HERE

Second, create a dataset call `d` by joining `d_return` and `d_tally` by column `name_image`. `d` contains rows whose identifiers appear in *both* datasets and columns from *both* datasets.

```
## # A tibble: 53,289 x 96
##   name_image state district fraud_proba fraud_bin foto seccion casilla dtto
##   <chr>      <chr> <chr>      <dbl> <lgl>      <chr> <chr> <chr> <chr>
## 1 aguascalien~ Agua~ I          8.04e-4 FALSE 2014~ 1      84 <NA>
## 2 aguascalien~ Agua~ I          4.28e-2 FALSE 2014~ 85     85 1
## 3 aguascalien~ Agua~ I          4.23e-1 FALSE 2014~ 45    45-A 1
## 4 aguascalien~ Agua~ I          3.49e-2 FALSE 2014~ 86     86 1
## 5 aguascalien~ Agua~ I          1.30e-1 FALSE 2014~ 87     87 1
## 6 aguascalien~ Agua~ I          2.12e-1 FALSE 2014~ 1     87-A 7
## 7 aguascalien~ Agua~ I          3.51e-2 FALSE 2014~ 88     88 1
## 8 aguascalien~ Agua~ I          3.19e-1 FALSE 2014~ 89     89 1
## 9 aguascalien~ Agua~ I          6.00e-8 FALSE 2014~ 89    89-A 7
## 10 aguascalien~ Agua~ I          3.60e-1 FALSE 2014~ 89    89-B 7
## # i 53,279 more rows
## # i 87 more variables: dto <roman>, municipio <chr>, edo <chr>, entidad <chr>,
## # pagina <dbl>, p1 <dbl>, p2 <dbl>, p3 <dbl>, p4 <dbl>, p5 <dbl>, pan <dbl>,
## # pri <dbl>, pps <dbl>, psm <dbl>, pms <dbl>, pfcrrn <dbl>, prt <dbl>,
## # parm <dbl>, noregis <dbl>, nombrenore <chr>, otros <dbl>, otroscan <chr>,
## # pan2 <dbl>, pri2 <dbl>, pps2 <dbl>, psm2 <dbl>, pms2 <dbl>, pfcrrn2 <dbl>,
## # prt2 <dbl>, parm2 <dbl>, noregis2 <dbl>, otro2 <dbl>, pan3 <dbl>, ...
```

Task 4. Visualize distributions of fraudulent tallies across candidates (6pt)

In this task, you will visualize the distributions of fraudulent tally sheets across three presidential candidates: **Sarinas (PRI)**, **Cardenas (FDN)**, and **Clouthier (PAN)**. The desired output of is reproducing and extending Figure 4 in the research article (Cantu 2019, pp. 720).

Task 4.1. Calculate vote proportions of Salinas, Clouthier, and Cardenas

Before getting to the visualization, you should first calculate the proportion of votes (among all) received by the three candidates of interest. As additional background information, there are two more presidential candidates in this election, whose votes received are recorded in `ibarra` and `castillo` respectively. Please perform the tasks in the following two steps on the `d` dataset:

- Create a new column named `total_president` as an indicator of the total number of votes of the 5 presidential candidates.
- Create three columns `salinas_prop`, `cardenas_prop`, and `clouthier_prop` that indicate the proportions of the votes these three candidates receive respectively.

Task 4.2. Replicate Figure 4

Based on all the previous step, reproduce Figure 4 in Cantu (2019, pp. 720).

Note: Your performance in this task will be mainly evaluated based on your output's similarity with the original figure. Pay attention to the details. For your reference, below is a version created by the instructor.

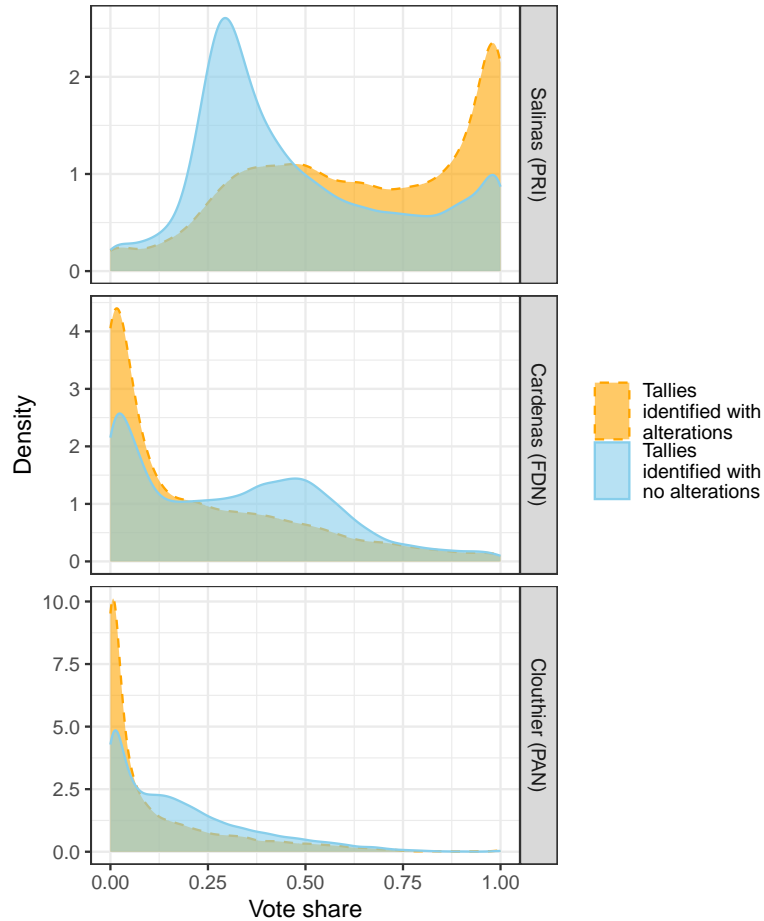


Figure 2: Distribution of Vote Shares for Each of the Candidates. Mexico, 1988

Task 4.3. Discuss and extend the reproduced figure

Referring to your reproduced figures and the research articles, in what way is the researcher's argument supported by this figure? Make an alternative visualization design that can substantiate and even augment the current argument. After you have shown your alternative design, in a few sentences, describe how your design provides visual aid as effectively as or more effectively than the original figure.

Note: Feel free to suggest *multiple* alternative designs to earn bonus credits. However, please be selective. Only a design with major differences from the existing ones can be counted as an alternative design.

YOUR CODE HERE

Task 5. Visualize the discrepancies between presidential and legislative Votes (6pt)

In this task, you will visualize the differences between the number of presidential votes across tallies. The desired output of is reproducing and extending Figure 5 in the research article (Cantu 2019, pp. 720).

Task 5.1. Get district-level discrepancies and fraud data

As you might have noticed in the caption of Figure 5 in Cantu (2019, pp. 720), the visualized data are aggregated to the *district* level. In contrast, the unit of analysis in the dataset we are working with, *d*, is *tally*. As a result, the first step of this task is to aggregate the data. Specifically, please aggregate *d* into a new data frame named `sum_fraud_by_district`, which contains the following columns:

- `state`: Names of states
- `district`: Names of districts
- `vote_president`: Total numbers of presidential votes
- `vote_legislature`: Total numbers of legislative votes
- `vote_diff`: Total number of presidential votes minus total number of legislative votes
- `prop_fraud`: Proportions of fraudulent tallies (hint: using `fraud_bin`)

Print your output data frame. For your reference, it should look like the following:

```
## # A tibble: 300 x 6
## # Groups:   state [32]
##   state      district vote_president vote_legislature vote_diff prop_fraud
##   <chr>      <chr>          <dbl>          <dbl>          <dbl>    <dbl>
## 1 Aguascalientes I          118139          102213          15926     0.135
## 2 Aguascalientes II         58722           55271           3451     0.215
## 3 Baja California I          75385           60550           14835     0.171
## 4 Baja California II         44630           32429           12201     0.0960
## 5 Baja California III        79072           75940           3132     0.132
## 6 Baja California IV       104627           90270           14357     0.375
## 7 Baja California V         55792           48971           6821     0.152
## 8 Baja California VI        64986           60596           4390     0.368
## 9 Baja Californi~ I         52226           47569           4657     0.259
## 10 Baja Californi~ II        30405           26641           3764     0.0933
## # i 290 more rows
```

Task 5.2. Replicate Figure 5

Based on all the previous step, reproduce Figure 5 in Cantu (2019, pp. 720).

Note 1: Your performance in this task will be mainly evaluated based on your output's similarity with the original figure. Pay attention to the details. For your reference, below is a version created by the instructor.

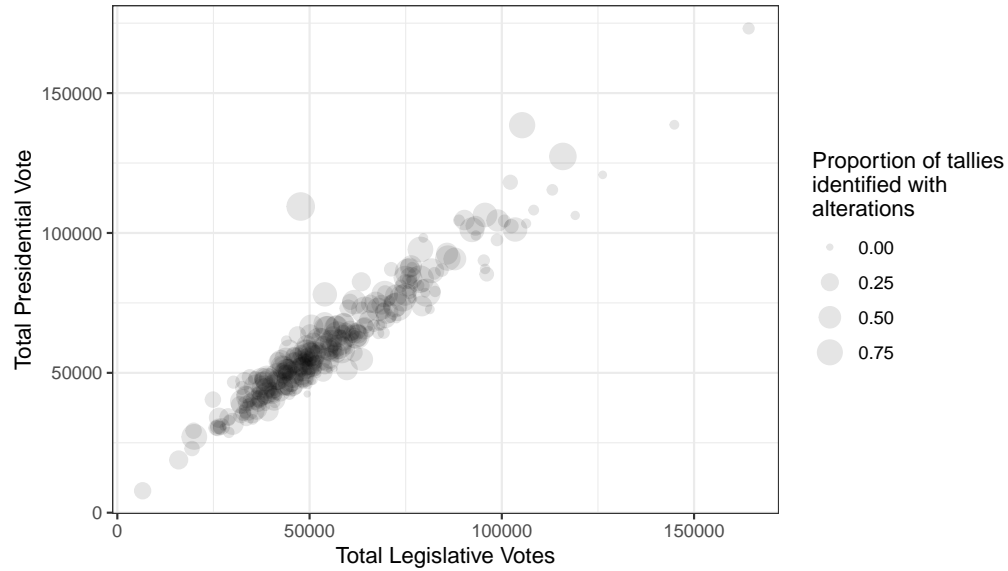


Figure 3: Total Number of District Votes for Presidential and Legislative Elections. Mexico, 1988

Note 2: The instructor has detected some differences between the above figure with Figure 5 on the published article. Please use the instructor's version as your main benchmark.

Task 5.3. Discuss and extend the reproduced figure

Referring to your reproduced figures and the research articles, in what way is the researcher's argument supported by this figure? Make an alternative visualization design that can substantiate and even augment the current argument. After you have shown your alternative design, in a few sentences, describe how your design provides visual aid as effectively as or more effectively than the original figure.

Note: Feel free to suggest *multiple* alternative designs to earn bonus credits. However, please be selective. Only a design with major differences from the existing ones can be counted as an alternative design.

YOUR CODE HERE

Task 6. Visualize the spatial distribution of fraud (6pt)

In this final task, you will visualize the spatial distribution of electoral fraud in Mexico. The desired output of is reproducing and extending Figure 3 in the research article (Cantu 2019, pp. 720).

Note 3. Load map data

As you may recall, map data can be stored and shared in **two** ways. The simpler format is a table where each row has information of a point that “carves” the boundary of a geographic unit (a Mexican state in our case). In this type of map data, a geographic unit is represented by multiple rows. Alternatively, a map can be represented by a more complicated and more powerful format, where each geographic unit (a Mexican state in our case) is represented by an element of a **geometry** column. For this task, I provide you with a state-level map of Mexico represented by both formats respectively.

Below the instructor provide you with the code to load the maps stored under the two formats respectively. Please run them before starting to work on your task.

```
# Load map (simple)
map_mex <- read_csv("data/map_mexico/map_mexico.csv")
# Load map (sf): You need to install and load library "sf" in advance
map_mex_sf <- st_read("data/map_mexico/shapefile/gadm36_MEX_1.shp")
map_mex_sf <- st_simplify(map_mex_sf, dTolerance = 100)
```

Bonus question: Explain the operations on `map_mex_sf` in the instructor’s code above.

Note: The map (sf) data we use are from https://gadm.org/download_country_v3.html.

Task 6.1. Reproduce Figure 3 with map_mex

In this task, you are required to reproduce Figure 3 with the `map_mex` data.

Note:

- Your performance in this task will be mainly evaluated based on your output's similarity with the original figure. Pay attention to the details. For your reference, below is a version created by the instructor.
- Hint: Check the states' names in the map data and the electoral fraud data. Recode them if necessary.

```
# YOUR CODE HERE
```

For your reference, below is a reproduced version created by the instructor.



Figure 4: Rates of Tallies Classified as Altered by State

Task 6.2. Reproduce Figure 3 with `map_mex_sf`

In this task, you are required to reproduce Figure 3 with the `map_mex` data.

Note:

- Your performance in this task will be mainly evaluated based on your output's similarity with the original figure. Pay attention to the details. For your reference, below is a version created by the instructor.
- Hint: Check the states' names in the map data and the electoral fraud data. Recode them if necessary.

YOUR CODE HERE

For your reference, below is a reproduced version created by the instructor.

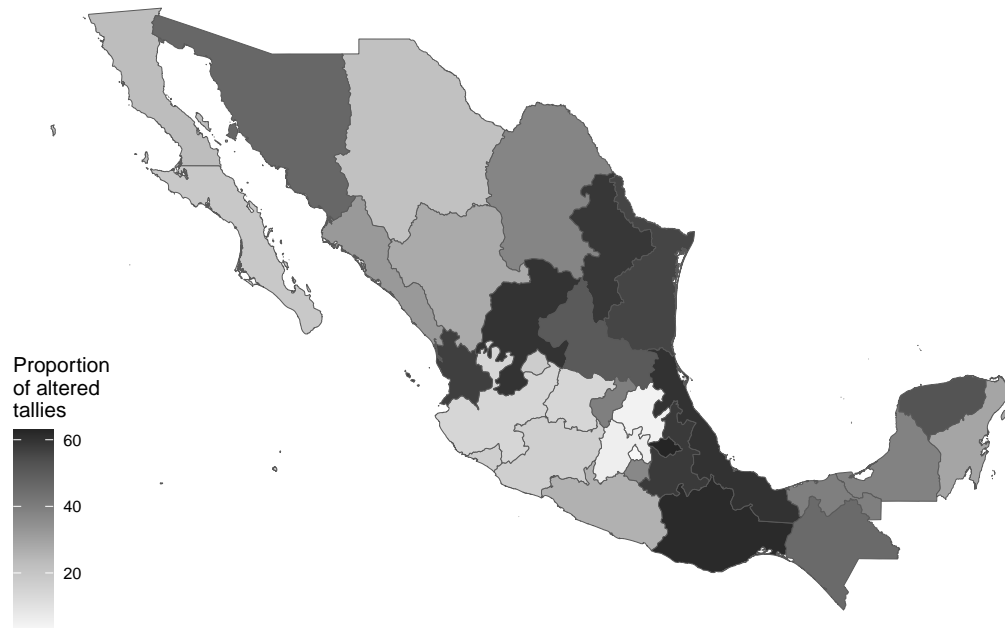


Figure 5: Rates of Tallies Classified as Altered by State

Task 6.3. Discuss and extend the reproduced figures

Referring to your reproduced figures and the research articles, in what way is the researcher's argument supported by this figure? Make an alternative visualization design that can substantiate and even augment the current argument. After you have shown your alternative design, in a few sentences, describe how your design provides visual aid as effectively as or more effectively than the original figure.

Note: Feel free to suggest *multiple* alternative designs to earn bonus credits. However, please be selective. Only a design with major differences from the existing ones can be counted as an alternative design.

YOUR CODE HERE