# Extracting Data in PDFs using R

Nathaniel Bechhofer

# Some Background on R

- R is a programming language designed for statistical computing
- Free, open source, large community
  - RStudio, a powerful IDE
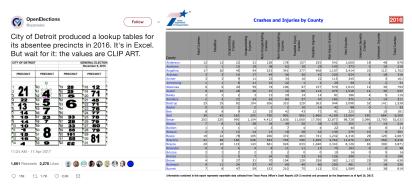  - CRAN (The Comprehensive R Archive Network)

# Examples of R packages

- `dplyr`, for a readable grammar of data manipulation
- `foreign`, for reading almost any kind of data file (including Stata files)
- `stringr`, for operating on strings
- `lubridate`, for parsing and processing dates
- `forecast`, for time series analysis
- `glmnet`, for regularized GLMs

along with over 10,000 other packages on CRAN.

# Using the `tabulizer` package to extract data from PDFs

- Many institutions (especially government agencies) publish data that may be useful for research

- Generally, ease of processing for end-users is not a priority

# Using the `tabulizer` package to extract data from PDFs

The `tabulizer` package for R allows you to solve some of these problems. (The Natural Resource Governance Institute has an online version with much of the functionality.)

Use case: you have data stuck in PDF tables

```r
library(tidyverse)
library(tabulizer)
tables_as_matrices <- tabulizer::extract_tables(filename_as_a_string,
                                    guess = TRUE, # restricts to capturing what it thinks is a table
                                    pages = c(1, 2, 4) # only searches pages 1, 2, and 4
                                    )
combined_matrix <- do.call(rbind, tables_as_matrices) # combines tables into one matrix
combined_df <- as_data_frame(combined_matrix)
combined_df[combined_df == ""] <- NA
combined_df <- na.omit(combined_df)
```

# The result

| 3 | Angelina | 17 | 20 | 45 | 54 | 233 | 341 | 277 | 408 | 1,107 | 3,414 | 23 | 112 | 1,702 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 4 | Aransas | 3 | 3 | 14 | 17 | 45 | 56 | 30 | 42 | 220 | 624 | 6 | 18 | 318 |
| 5 | Archer | 3 | 3 | 9 | 11 | 23 | 34 | 10 | 22 | 115 | 245 | 1 | 3 | 161 |
| 6 | Armstrong | 1 | 1 | 8 | 11 | 10 | 22 | 1 | 2 | 29 | 69 | 2 | 2 | 51 |
| 7 | Atascosa | 6 | 6 | 40 | 59 | 74 | 106 | 97 | 157 | 570 | 1,613 | 12 | 30 | 799 |
| 8 | Austin | 9 | 10 | 28 | 36 | 47 | 72 | 66 | 113 | 475 | 1,518 | 12 | 36 | 637 |
| 9 | Bailey | 2 | 2 | 2 | 3 | 13 | 20 | 14 | 14 | 81 | 186 | 4 | 5 | 116 |
| 10 | Bandera | 2 | 2 | 23 | 24 | 62 | 84 | 37 | 57 | 183 | 375 | 4 | 14 | 311 |
| 11 | Bastrop | 25 | 35 | 82 | 104 | 206 | 310 | 229 | 363 | 944 | 3,098 | 52 | 141 | 1,538 |
| 12 | Baylor | 0 | 0 | 2 | 2 | 1 | 2 | 10 | 16 | 42 | 88 | 0 | 1 | 55 |
| 13 | Bee | 8 | 8 | 9 | 12 | 25 | 42 | 43 | 72 | 91 | 225 | 5 | 10 | 181 |
| 14 | Bell | 39 | 42 | 162 | 205 | 730 | 991 | 965 | 1,466 | 4,109 | 13,564 | 195 | 664 | 6,200 |
| 15 | Bexar | 203 | 220 | 995 | 1,154 | 4,412 | 5,836 | 11,060 | 17,766 | 32,873 | 98,728 | 3,090 | 13,700 | 52,633 |
| 16 | Blanco | 7 | 8 | 16 | 28 | 34 | 49 | 20 | 38 | 133 | 355 | 6 | 14 | 216 |
| 17 | Borden | 1 | 1 | 1 | 1 | 3 | 7 | 4 | 4 | 22 | 40 | 0 | 0 | 31 |
| 18 | Bosque | 2 | 2 | 12 | 14 | 13 | 18 | 28 | 42 | 136 | 279 | 10 | 8 | 201 |
| 19 | Bowie | 19 | 22 | 78 | 105 | 266 | 372 | 463 | 741 | 1,232 | 4,316 | 29 | 129 | 2,087 |