# Digitizing Maine's Voting History with a Statistical Analysis of Error Rate

By Gillian King
Bowdoin College '22
Kufe Family Research Fellowship

# Digitizing Raw Data

- **Goals:**
  - **Create an online database that stores as much digitized, publicly accessible data as possible in the state of Maine by using an Optical Character Recognition (OCR) software.**
  - **Analyze the efficiency of three statistical tests used to correct errors of the OCR process.**
    - Digital Maine Repository (handwritten data from 1984 to 1860s and before)
      - Digitization process varies state by state.
    - ABBYY Finereader PDF 15 OCR software to convert handwritten documents to searchable PDFs.

1980 General Election Results for Minnesota (Growe 1980, 16).

**Here, each column header represents yes/no vote totals for each question, and each row represents one town within a given county.**

1980 General Election Results for Maine (Bureau of Corporations, Elections and Commissions 1980).

# Analyzing the Error Rates of Data using OCR

I. Calculating the p-values of the data
   ● Examining patterns down a column of data (examining the totals by county for each question)
II. Summing the columns of the scanned files
   ● Looking for correct vote totals in the scanned files (comparing to raw data totals)
III. Random Spot-Check on R-Studio
   ● Running the code to produce one random column and row in scanned data to compare to raw sheet.

# What is a p-value?

- In hypothesis testing, there exists a **null hypothesis** and an **alternative hypothesis.**
    - In general, the null hypothesis suggests that an outcome of an experiment is random, while the alternative hypothesis suggests otherwise.
    - "Informally, a *p*-value is the probability under a specified statistical model that a statistical summary of the data (e.g., the sample mean difference between two compared groups) would be equal to or more extreme than its observed value" (Wasserstein 2016, 131).
    - Further discussion on its potential limitations later in presentation.

# P–Values in this presentation

- Goal: Is there a correlation between a county's proportion of results for one question when compared to their results for another?
  - How likely is the same county likely to vote yes/no to one question given that they voted yes/no to another?
- **In this presentation, the null hypothesis suggests that the entries after the OCR are correct, and that the proportion of yes/no votes per county stays consistent despite a difference in turnout over time.**
- **The alternative hypothesis suggests that for any individual entry after the OCR process, that entry is likely to be an error.**
- Thus, the lower the p-value for a row of data, the more likely that an entry was an error.

# I. Calculate the p–values of the data

Approach using R-Studio:
I. Created a list of yes/no vote totals for on entire question across all all counties for two questions at at time (in log form).
II. Plotted questions against each other to get a general sense of the strength of correlation.
III. Used a simple linear regression to create a line of best fit.
IV. Calculated p-values for each column of data using pnorm in RStudio.
V. Based on calculations in Part I, determined whether or not to re-analyze certain data points.

# Calculating the Least Squares Estimate (LSE)

- Use the LSE to create a line of best-fit for the data.
- Here, this is merely used to get a sense of large outliers in the data.
- LSE attempts to minimize "the sum of the squares of the errors that are generated by the results of the associated equations, such as the squared residuals resulting from differences in the observed value, and the value anticipated, based on that model" (Kelton and Kindness 2020).

```
a <- cov(Q.4,Q.5,use='complete')/var(Q.4,use='complete')

b <- mean(Q.5,na.rm=TRUE)-a*mean(Q.4,na.rm=TRUE)

sigma.squared <- (1/N)*sum((Q.5-a*Q.4-b)^2,na.rm=TRUE)
s <- sqrt(sigma.squared)
abline(b,a)

y.pred <- a*Q.4 + b

p.values <- pnorm(-abs(y.pred-Q.5),0,s)*2
plot(Q.4,log10(p.values),ylim=c(-15,0))
```

# Sample linear regression, data from Hancock County Sheets (1978)



Question 5 vs. Question 4

# Graphical Interpretation of p–values



Graphical Interpretation of p-values

# Summing the Columns

## Random Spot-Check in RStudio



(Figure 1)

(Bureau of Corporations, Elections and Commissions, 1978)

| | YES |
|---|---|
| Swan's Island, | 92 |
| Tremont, | 164 |
| Trenton, | 132 |
| Verona, | 93 |
| Waltham, | 34 |
| Winter Harbor, | 96 |
| PLANTATIONS | |
| Great Pond, | 18 |
| Long Island, | 10 |
| | 7697 |

(Figure 2)

| | YES.1 |
|---|---|
| Amherst | 40 |
| Aurora | 27 |
| Bar Harbo | 924 |
| Blue Hill | 352 |
| Brooklin | 127 |
| Brooksvil | 163 |
| Bucksport | 650 |
| Castine | 204 |
| Cranberry | 63 |
| Dedham | 166 |
| Deer Isle | 281 |
| Eastbrook | 40 |
| Ellsworth | 965 |
| Franklin | 122 |
| Gouldsbo | 190 |
| Hancock | 198 |
| Lamoine | 202 |
| Mariavill | 18 |
| Mount De | 460 |
| Orland | 260 |
| Osborn | 11 |
| Otis | 39 |
| Penobsco | 151 |
| Sedgwick | 137 |
| Sorrento | 66 |
| Southwes | 294 |
| Stoningto | 184 |
| Sullivan | 143 |
| Surry | 173 |
| Swan's Isl | 80 |
| Tremont | 147 |
| Trenton | 131 |
| Verona | 80 |
| Waltham | 18 |
| Winter Ha | 89 |
| Great Pon | 14 |
| Long Islar | 4 |

```
data.set <- read.table("1984.Androscoggin.csv",header=TRUE,sep=",")

num.row <- dim(data.set)[1]
num.col <- dim(data.set)[2]

random.col <- sample(2:num.col,1)
random.row <- sample(1:num.row,1)

print(as.character(data.set[random.row,1]))
print(as.numeric(as.character(data.set[random.row,random.col])))
print(c(random.row,random.col))
```

(Figure 3)

```
> source('C:/Users/gilli/Dropbox/1984 Documents copy/1984.Aroostook/SPOT CHECK.R')
[1] "Mariaville "
[1] 34
[1] 18  5
> source('C:/Users/gilli/Dropbox/1984 Documents copy/1984.Aroostook/SPOT CHECK.R')
[1] "Stonington "
[1] 196
[1] 27  3
> source('C:/Users/gilli/Dropbox/1984 Documents copy/1984.Aroostook/SPOT CHECK.R')
[1] "Deer Isle"
[1] 281
[1] 11  2
> source('C:/Users/gilli/Dropbox/1984 Documents copy/1984.Aroostook/SPOT CHECK.R')
[1] "Surry "
[1] 204
[1] 29 14
> |
```

(Figure 4)

# Discussions Around the Limitations of p-Values

Per the American Statistical Association (ASA):
- *P*-values do not directly "measure the probability that the studies hypothesis is true, or the probability that the data were produced by random chance alone" (Wasserstein and Lazar 2020, 131).
- Arbitrary cut-off point
  - Limitations in statistical studies of human behavior such as voting patterns (Wasserstein and Lazar 2020, 131).
- Effect size can provide a false sense of deviation from the null hypothesis (Wasserstein and Lazar 2020, 132).

# General Results and Research Design Improvements

**General Results:**
- Time spent: Test #2 > Test #1 > Test #3
- Errors Caught: Test #2 > Test #1 > Test #3

**Research Design Improvements:**
- **Crowd-sourcing**
  - Weighing the efficiency of voting results.
  - Will crowdsourcing be more efficient than using the OCR Software?
- Working with Maine State Archives to order printed voting data.

# Bibliography

Bureau of Corporations, Elections and Commissions. "1980 General Election: Presidential," 1980.

Growe, Joan Anderson. *Minnesota Election Results 1980: Primary Election and General Election, Presidential Electors, Representatives in Congress, Senators in Legislature, Representatives in Legislature, Supreme Court, Constitutional Amendments*. St. Paul: Office of the Secretary of State, Election Division, 1980.

Kenton, Will, and David Kindness. "How the Least Squares Method Works." Investopedia. Investopedia, September 16, 2020. https://www.investopedia.com/terms/l/least-squares-method.asp.

Wasserstein, Ronald L. "ASA Statement on Statistical Significance and P-Values." *The American Statistician* 70, no. 2 (2016).

Wasserstein, Ronald L., and Nicole A. Lazar. "The ASA Statement on p-Values: Context, Process, and Purpose." *The American Statistician* 70, no. 2 (2016): 129–33. https://doi.org/10.1080/00031305.2016.1154108.