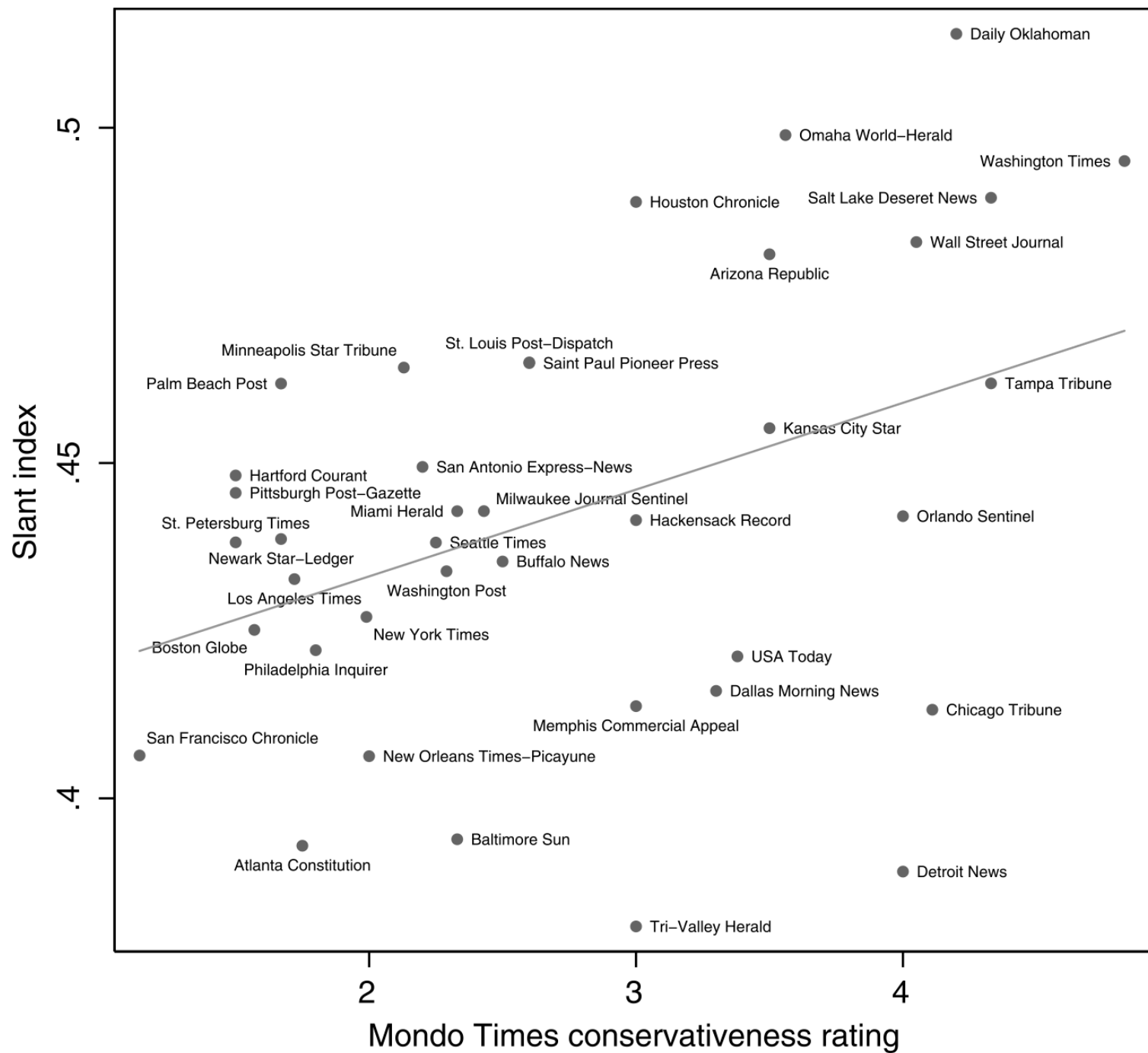# Measuring bias in Indian newspapers

Arun Varghese

# Inspiration

- Gentzkow & Shapiro 2010
- First economics paper to use text data in an innovative way
- 'What drives media slant?': supply side or demand side?
- Intermediate step: need to quantify 'slant' or bias
- Done by comparing text from congressional record (ideology labeled using constituency vote share) against newspaper text

Daily Oklahoman

Omaha World–Herald

Washington Times

Houston Chronicle

Salt Lake Deseret News

Wall Street Journal

Arizona Republic

Minneapolis Star Tribune

St. Louis Post–Dispatch

Saint Paul Pioneer Press

Palm Beach Post

Tampa Tribune

Kansas City Star

Hartford Courant

San Antonio Express–News

Pittsburgh Post–Gazette

Milwaukee Journal Sentinel

Miami Herald

Orlando Sentinel

Hackensack Record

St. Petersburg Times

Newark Star–Ledger

Seattle Times

Buffalo News

Los Angeles Times

Washington Post

New York Times

Boston Globe

USA Today

Philadelphia Inquirer

Dallas Morning News

Memphis Commercial Appeal

Chicago Tribune

San Francisco Chronicle

New Orleans Times–Picayune

Baltimore Sun

Atlanta Constitution

Detroit News

Tri–Valley Herald

Slant index

.5

.45

.4

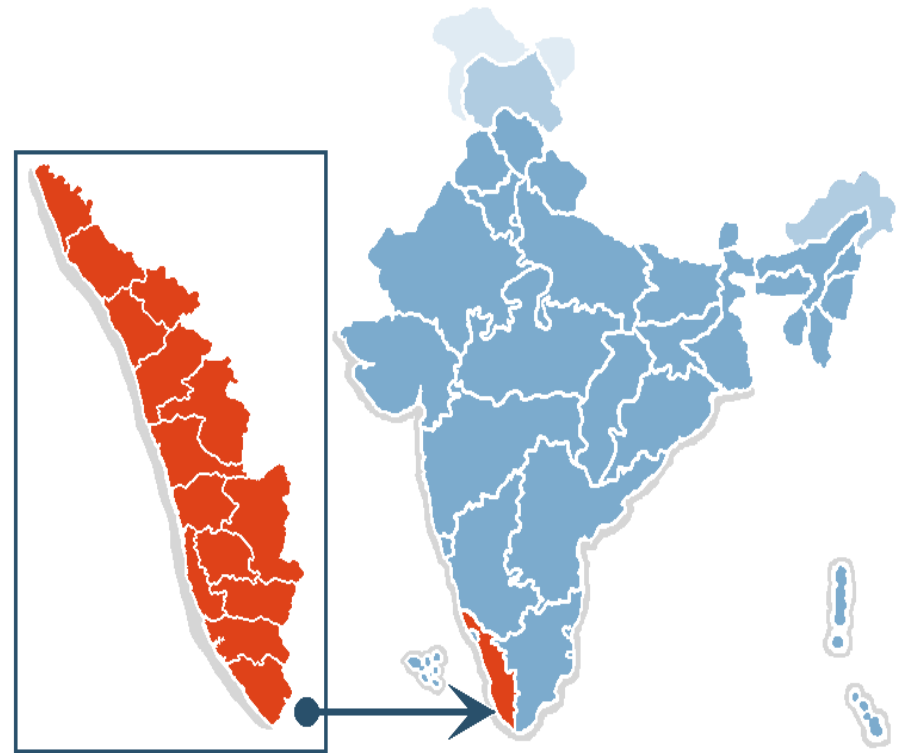Mondo Times conservativeness rating

2

3

4

# My project

- Similar exercise with Malayalam-language newspapers in the Indian state of Kerala

- Want to get better at language, learn something about politics

- Challenges: dealing with non-standard text, not many pre-existing NLP resources, agglutinative language

# Kerala

- South-western tip of India
- First state to democratically elect communist government
- Now, multiple parties but dominated by two coalitions: Left Democratic Front (LDF) and congress-led United Democratic Front (UDF) who have traded majority power in legislature since Kerala founding

# More background

- Legislative Assembly: Niyamasabha (140 elected members)
- Last election in 2011
- Print media widely consumed
- Bias is widely thought to exist – some papers are owned by certain parties

# Method

- Get text from legislative hearings
- Label by political slant/ideology from vote share data
- Generate text features and build regression model of slant
- Get newspaper article text and fit to model to get slant score for each newspaper

# Data

- Text of legislative meetings available online
  - But mostly as PDFs
  - But 'Question and Answer' sessions available on html pages (bad sample?)
  - Scraped a month's worth of these (labeled with speaker/s)
- Vote share data from 2011 election available online
- Text from newspaper articles available online

# Feature engineering and selection

- Represent legislative text as bag of words
- Drop 'stop words', rare words (have to define own stop words)
- Each document has ideology label
- Use chi-squared test to find words/phrases most representative of a particular party

$$\chi^2 = \sum \frac{(Observed - Expected)^2}{Expected}$$

- Intuition: If we assumed party and words used were randomly assigned, which co-occur more/less often than we would expect?

# Gentzkow & Shapiro

**Democrats:**

- 'Workers rights'
- 'Poor people'
- 'Estate tax'

**Republicans:**

- 'Stem cell'
- 'Saddam Hussein'
- 'Death tax'

# Gentzkow & Shapiro

# Me

**Democrats:**

- 'Workers rights'
- 'Poor people'
- 'Estate tax'

**Republicans:**

- 'Stem cell'
- 'Saddam Hussein'
- 'Death tax'

Lots of procedural words, question words, etc

**LDF:**

- അഗ്രികുൽറ്റു (agriculture)
- ആദിവാസി (adivasi)
- More negative words [അവബോധ, അനധികൃത]

**UDF:**

- Many procedural/question words
- അഞ്ചുവർഷങ്ങളിലായി വകയിരുത്തി (For five years, we've been here)
- And variations

# Ongoing work

- Tinkering with model – not getting great explanatory power in-sample

- Use article text to help generate features

- Fit article text and get slant of newspapers!

# Thanks!