# Applications: Prediction

Matthew Gentzkow
Jesse M. Shapiro

Chicago Booth and NBER

# Introduction

- Methods map high-dimensional $x$ to low-dimensional $z$
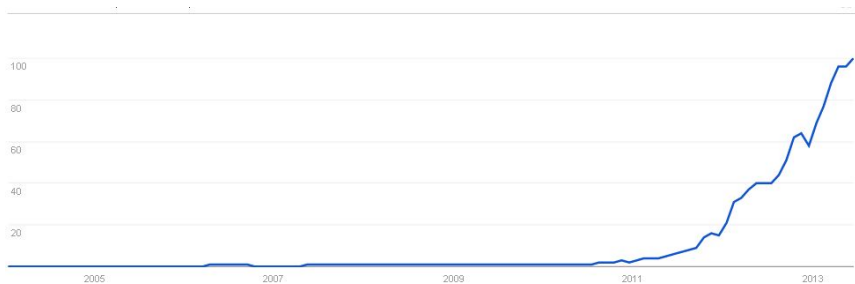
# Introduction

- Methods map high-dimensional $x$ to low-dimensional $z$
- Three main uses of $z$
  1. Forecasting (e.g., what will inflation be next month?)
  2. Descriptive analysis (e.g., are there genes that predict risk aversion?)
  3. Input into subsequent causal analysis (e.g., as LHS var, RHS var, control, instrument, etc.)

# Outline

- Brief overview of data & applications
- Detailed discussion of text as data

# Overview

"Big Data"

# Google Searches



"Big Data"

- Prediction
    - Google flu trends (Dukik et al. 2009)
    - Unemployment claims, retail sales, consumer confidence, etc. (Choi & Varian 2009, 2012)

# Google Searches: Applications

- Prediction
  - Google flu trends (Dukik et al. 2009)
  - Unemployment claims, retail sales, consumer confidence, etc. (Choi & Varian 2009, 2012)

- Descriptive
  - What searches predict "consumer confidence" (Varian 2013)

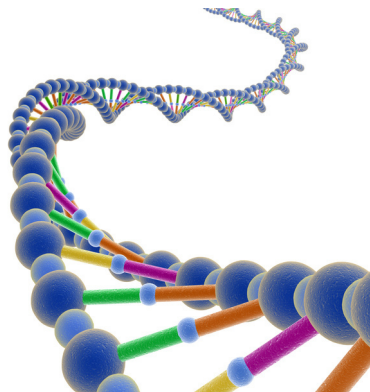# Google Searches: Applications

- Prediction
  - Google flu trends (Dukik et al. 2009)
  - Unemployment claims, retail sales, consumer confidence, etc. (Choi & Varian 2009, 2012)

- Descriptive
  - What searches predict "consumer confidence" (Varian 2013)

- Input to analysis
  - Saiz & Simonsohn (2013) $\rightarrow$ city-level corruption
  - Stephens-Davidowitz (2013) $\rightarrow$ racial animus & effect on voting for Obama

- Genetic data has been one of the main applications of high-dimensional methods
- LHS: Physical or behavioral outcome
- RHS: Single-nucleotide polymorphisms (SNPs)
- Typical dataset is $N \approx 10,000$ and $K \approx 2,500,000$

# Genes: Applications

- Descriptive: look for genetic predictors of...
  - Risk aversion & social preferences (Cesarini et al. 2009)
  - Financial decision making (Cesarini et al. 2010)
  - Political preferences (Benjamin et al. 2012)
  - Self-employment (van der Loos et al. 2013)
  - Educational attainment, subjective well being (Rietveld et al. forthcoming)

# Genes: Applications

- Descriptive: look for genetic predictors of...
  - Risk aversion & social preferences (Cesarini et al. 2009)
  - Financial decision making (Cesarini et al. 2010)
  - Political preferences (Benjamin et al. 2012)
  - Self-employment (van der Loos et al. 2013)
  - Educational attainment, subjective well being (Rietveld et al. forthcoming)
- Early reported associations have been shown to be spurious & non-replicable (Benjamin et al. 2012)

- Big and high dimensional
  - 10 years of Medicare data on the order of 100 TB
  - (*patient* × *doctor* × *hospital* × *treatment* × *cost*...)

# Medical Claims



- Big and high dimensional
  - 10 years of Medicare data on the order of 100 TB
  - (*patient* × *doctor* × *hospital* × *treatment* × *cost*...)
- Dimension reduction: How to collapse data into a single-dimensional index of "health" or "predicted spending"
  - Medicare "risk scores" based on ad hoc criteria
  - Johns Hopkins ACG system uses proprietary predictive model

# Medical Claims: Applications



- Input into analysis
    - Numerous studies use Medicare risk scores as a control variable or independent variable of interest
    - Einav & Finkelstein (forthcoming) use risk score as mediator of health plan choice
    - Handel (2013) uses Johns Hopkins ACG score as measure of private information

# Credit Scores



- Predicting default risk from consumer credit data is similar to predicting health risk from medical claims

# Credit Scores



- Predicting default risk from consumer credit data is similar to predicting health risk from medical claims
- Forecasting
  - Large literature applies machine learning tools to improve forecasting of default risks (e.g., Khandani et al. 2010)

# Credit Scores



- Predicting default risk from consumer credit data is similar to predicting health risk from medical claims
- Forecasting
  - Large literature applies machine learning tools to improve forecasting of default risks (e.g., Khandani et al. 2010)
- Input into analysis
  - Adams, Einav and Levin (2009) and Einav, Jenkins and Levin (2012) evaluate auto dealer's proprietary credit scoring algorithm
  - Rajan, Seru and Vig (forthcoming) look at market responses to using a limited set of variables in credit scoring

- Amazon, Ebay, and other large Internet firms use purchase and browsing history to make recommendations, target advertising, etc.
- "Netflix Prize" contest for best algorithm to predict future user ratings based on past ratings

- Poole & Rosenthal (1984, 1985, 1991, 2000, etc.) use factor analysis methods to project Roll Call votes into ideology scores
- Ask questions like
  - Are there multiple dimensions of ideology?
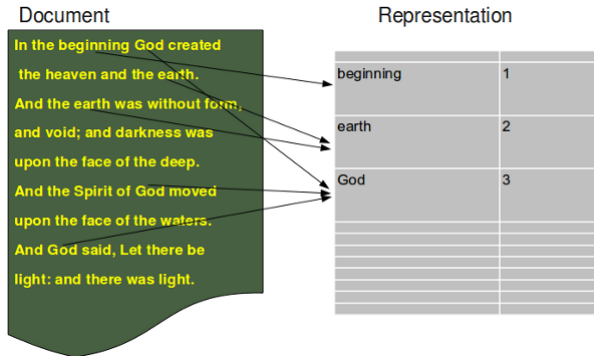  - How has polarization changed over time?

# Text as Data

# Sources

- News
- Books
- Web content
- Congressional speeches
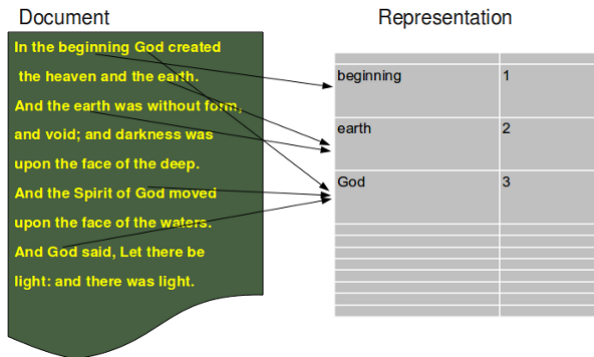- Corporate filings
- Twitter & Facebook

- Amazon and eBay listings
- Google search ads
- Medical records
- Central bank announcements
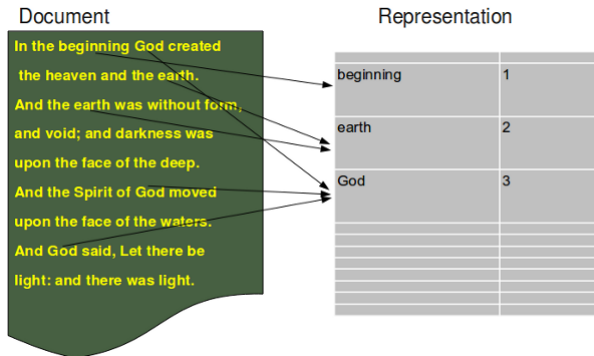
# Bag of Words

# Bag of Words



- Can apply to "*N*-grams" as well as single words

# Bag of Words



- Can apply to "*N*-grams" as well as single words
- This seems crude, but it works remarkably well in practice, and gains to more sophisticated representations prove to be small

- General theme: Real research always combines automated dimension reduction techniques with "manual" steps based on priors

# The Science and the Art

- General theme: Real research always combines automated dimension reduction techniques with "manual" steps based on priors
- E.g.,
  - Keep only words occurring more than $X$ times
  - Drop very common "stopwords" like "the," "at," "a"
  - "Stem" words to combine, e.g., "economics," "economic," "economically"
  - Drop HTML tags

- General theme: Real research always combines automated dimension reduction techniques with "manual" steps based on priors
- E.g.,
  - Keep only words occurring more than $X$ times
  - Drop very common "stopwords" like "the," "at," "a"
  - "Stem" words to combine, e.g., "economics," "economic," "economically"
  - Drop HTML tags
- In fact, 90% of text analysis in economics does not automated dimension reduction at all
  - Saiz & Simonsohn (2013) $\rightarrow$ city name + "corruption"
  - Baker, Bloom, and Davis (2013) $\rightarrow$ "economic" + "policy" + "uncertainty", etc.
  - Lucca & Trebbi (2011) $\rightarrow$ "hawkish/dovish," "loose/tight," + "Federal Open Market Committee"

# Interfaces

- Bag of words assumes access to full text (at least for $N$-grams of interest)

# Interfaces

- Bag of words assumes access to full text (at least for *N*-grams of interest)
- Many researchers, however, can only access text via *search* interfaces (e.g., Google page counts, news archives, etc.)

# Interfaces

- Bag of words assumes access to full text (at least for *N*-grams of interest)
- Many researchers, however, can only access text via *search* interfaces (e.g., Google page counts, news archives, etc.)
  - This requires some external method of feature selection to narrow down vocabulary

# Sentiment Analysis

# Setup

- Outcome $y_i$
- Features $x_i$
- Data:

$$\underbrace{\{x_1, y_1\}, \{x_2, y_2\}, \{x_3, y_3\}, ..., \{x_N, y_N\}}_{\text{Training set}}, \underbrace{\{x_{N+1}, ?\}}_{\text{Target}}$$

# Spam Filter

- Outcome $y_i \in \{spam, ham\}$
- Human coder classifies $N$ cases as *spam* or *ham*
- Must decide whether to deliver the $(N + 1)$ message or send it to the filter

- What features $x_i$ do we use?
    - Counts of words?
    - Counts of characters?
    - Complete machine representation of e-mail?

- How do we avoid overfitting?
    - >1m words in English language
    - ASCII file with 100 printable characters has $95^{100}$ possible realizations

# Applications

- Partisanship in the news media [TODAY]
  - Turn millions of words into an index of media slant or bias
- Sentiment in financial news [TODAY]
  - Classify news, chat room discussions, etc. as positive or negative
- Estimating causal effects [TOMORROW]
  - Turn a huge dataset into a low-dimensional control for endogeneity

# Sentiment Analysis: *Partisanship in the News Media*

- Questions
  - How centrist are the news media?
  - What factors (owners, readers) predict how a newspaper portrays the news?
- Need measure of partisan orientation of news media
- Challenges
  - Training set: research assistants? surveys?
  - Dimensionality
    - Feature selection: words? phrases? images?
    - Parsimony: millions of possible words/phrases

- Training set: US Congress
  - Assign members an ideology score $y_i$ based on roll-call voting
- Dimensionality
  - Count frequency of citations to think tanks $x_i$ in Congressional Record
  - Feature selection "by hand": use *ex ante* criterion to control dimensionality

# Groseclose and Milyo (2005)

**Search the Congressional Record**
**105th Congress** (1997-1998)
The Congressional Record is the official record of the proceedings and debates of the U.S. Congress. ● More about the Congressional Record

Print    Subscribe    Share/Save

Search the Congressional Record | Latest Daily Digest | Browse Daily Issues
Browse the Keyword Index | Congressional Record App

**Select Congress:**
113 | 112 | 111 | 110 | 109 | 108 | 107 | 106 | 105 | 104 | 103 | 102 | 101

Congress-to-Year Conversion

● Help

**Enter Search**    SEARCH    CLEAR

brookings institution

○ Exact Match Only    ● Include Variants (plurals, etc.)

**Member of Congress**

| Any Representative | | Any Senator | |
|---|---|---|---|
| Abercrombie, Neil (HI-1) | | Abraham, Spencer (MI) | |
| Ackerman, Gary L. (NY-5) | | Akaka, Daniel K. (HI) | |
| Aderholt, Robert B. (AL-4) | | Allard, Wayne (CO) | |

○ Member speaking or mentioned    ● All occurrences

Mr. DORGAN. Mr. President, I come to the floor to speak first about the Congressional Budget Office, which last week released its monthly budget projection. And I noticed that this projection, this estimate, received prominent coverage in the Washington Post and in other major daily newspapers around the country last week....

A study by a tax expert at the **Brookings Institution** says if you have a national sales tax, the rates would probably be over 30 percent, and then add the State and local taxes, and that would be on almost everything. So say you would like to buy a house and here is the price we have agreed on, and then have someone tell you, oh, yes, you have a 37-percent sales tax applied to that price, 30 percent Federal, 7 percent State and local.

| Speaker | Think Tank | Count |
|---|---|---|
| Dorgan, Byron (D-ND) | Brookings Institution | 1 |

- Count references to think tanks in news media

- Let $y_i$ be ADA score of senator $i$
- Let $x_{ijt}$ be indicator for senator $i$ cites think tank $j$ on occasion $t$
- Then
$$\Pr\left(x_{ijt} = 1\right) = \frac{\exp\left(\alpha_j + \beta_j y_i\right)}{\sum_{j'} \exp\left(\alpha_{j'} + \beta_{j'} y_i\right)}$$
- Assume same model applies to news media $m$ but treat $y_m$ as unknown
- Estimate $\alpha_j, \beta_j, y_m$ via joint maximum likelihood

- Dimension of data $p = 50$
  - Start with 200 think tanks
  - Collapse all but top 44 into 6 groups
- Dimension of data $n = 535$
  - Less those who don't cite think tanks

# Groseclose and Milyo (2005)

| Senator | Partisanship | News Outlet | Partisanship |
|---------|--------------|-------------|--------------|
| John McCain | 12.7 | Fox News | |
| Arlen Specter | 51.3 | *USA Today* | |
| Joe Lieberman | 74.2 | *New York Times* | |

# Groseclose and Milyo (2005)

| Senator | Partisanship | News Outlet | Partisanship |
|---------|--------------|-------------|--------------|
| John McCain | 12.7 | Fox News | 39.7 |
| Arlen Specter | 51.3 | *USA Today* | 63.4 |
| Joe Lieberman | 74.2 | *New York Times* | 73.7 |

# Groseclose and Milyo (2005)

- Text of 2005 Congressional Record
- Scripted pipeline:
    - Download text
    - Split up text into individual speeches
    - Identify speaker
    - Count *all* two-word/three-word phrases

- Training set: US Congress
  - Assign members an ideology score $y_i$ based on partisanship of constituents

- Dimensionality
  - Compute frequency table of phrase counts by party
  - Compute $\chi^2$ statistic of independence
  - Identify 1000 phrases with highest $\chi^2$

- Memo to Rep. candidates: "Never say '**privatization/private accounts**.' Instead say '**personalization/personal accounts**.' Two-thirds of America want to personalize Social Security while only one-third would privatize it. Why? Personalizing Social Security suggests ownership and control over your retirement savings, while privatizing it suggests a profit motive and winners and losers."

# Example: Social Security

- Memo to Rep. candidates: "Never say '**privatization/private accounts**.' Instead say '**personalization/personal accounts**.' Two-thirds of America want to personalize Social Security while only one-third would privatize it. Why? Personalizing Social Security suggests ownership and control over your retirement savings, while privatizing it suggests a profit motive and winners and losers."

- Congress: "personal account" (48 D vs 184 R); "private account" (542 D vs 5 R)

# Top Phrases

| Republicans: 2-word | Republicans: 3-word | Democrats: 2-word | Democrats: 3-word |
|---|---|---|---|
| stem cell | embryonic stem cell | private accounts | veterans health care |
| natural gas | hate crimes legislation | trade agreement | congressional black caucus |
| death tax | adult stem cells | american people | va health care |
| illegal aliens | oil for food program | tax breaks | billion in tax cuts |
| class action | personal retirement accounts | trade deficit | credit card companies |
| war on terror | energy and natural resources | oil companies | security trust fund |
| embryonic stem | global war on terror | credit card | social security trust |
| tax relief | hate crimes law | nuclear option | privatize social security |
| illegal immigration | change hearts and minds | war in iraq | american free trade |
| date the time | global war on terrorism | middle class | central american free |
| boy scouts | class action fairness | african american | national wildlife refuge |
| hate crimes | committee on foreign relations | budget cuts | dependence on foreign oil |
| oil for food | deficit reduction bill | nuclear weapons | tax cuts for the wealthy |
| global war | boy scouts of america | checks and balances | vice president cheney |
| medical liability | repeal of the death tax | civil rights | arctic national wildlife |
| highway bill | highway trust fund | veterans health | bring our troops home |
| adult stem | action fairness act | cut medicaid | social security privatization |
| democratic leader | committee on commerce science | foreign oil | billion trade deficit |
| federal spending | cord blood stem | president plan | asian pacific american |
| tax increase | medical liability reform | gun violence | president bush took office |

# Top Phrases: Social Security

| Republicans: 2-word | Republicans: 3-word | Democrats: 2-word | Democrats: 3-word |
|---|---|---|---|
| stem cell | embryonic stem cell | **private accounts** | veterans health care |
| natural gas | hate crimes legislation | trade agreement | congressional black caucus |
| death tax | adult stem cells | american people | va health care |
| illegal aliens | oil for food program | tax breaks | billion in tax cuts |
| class action | **personal retirement accounts** | trade deficit | credit card companies |
| war on terror | energy and natural resources | oil companies | **security trust fund** |
| embryonic stem | global war on terror | credit card | **social security trust** |
| tax relief | hate crimes law | nuclear option | **privatize social security** |
| illegal immigration | change hearts and minds | war in iraq | american free trade |
| date the time | global war on terrorism | middle class | central american free |
| boy scouts | class action fairness | african american | national wildlife refuge |
| hate crimes | committee on foreign relations | budget cuts | dependence on foreign oil |
| oil for food | deficit reduction bill | nuclear weapons | tax cuts for the wealthy |
| global war | boy scouts of america | checks and balances | vice president cheney |
| medical liability | repeal of the death tax | civil rights | arctic national wildlife |
| highway bill | highway trust fund | veterans health | bring our troops home |
| adult stem | action fairness act | cut medicaid | **social security privatization** |
| democratic leader | committee on commerce science | foreign oil | billion trade deficit |
| federal spending | cord blood stem | president plan | asian pacific american |
| tax increase | medical liability reform | gun violence | president bush took office |

Other R: **personal accounts; social security reform; social security system**

Other D: **privatization plan; security trust; security trust fund; social security trust; privatize social security; social security privatization; privatization of social security; cut social security**

# Top Phrases: Foreign Policy

| Republicans: 2-word | Republicans: 3-word | Democrats: 2-word | Democrats: 3-word |
|---|---|---|---|
| stem cell | embryonic stem cell | private accounts | **veterans health care** |
| natural gas | hate crimes legislation | trade agreement | congressional black caucus |
| death tax | adult stem cells | american people | **va health care** |
| illegal aliens | **oil for food program** | tax breaks | billion in tax cuts |
| class action | personal retirement accounts | trade deficit | credit card companies |
| **war on terror** | energy and natural resources | oil companies | security trust fund |
| embryonic stem | **global war on terror** | credit card | social security trust |
| tax relief | hate crimes law | nuclear option | privatize social security |
| illegal immigration | **change hearts and minds** | **war in iraq** | american free trade |
| date the time | **global war on terrorism** | middle class | central american free |
| boy scouts | class action fairness | african american | national wildlife refuge |
| hate crimes | committee on foreign relations | budget cuts | **dependence on foreign oil** |
| **oil for food** | deficit reduction bill | **nuclear weapons** | tax cuts for the wealthy |
| **global war** | boy scouts of america | checks and balances | vice president cheney |
| medical liability | repeal of the death tax | civil rights | arctic national wildlife |
| highway bill | highway trust fund | veterans health | **bring our troops home** |
| adult stem | action fairness act | cut medicaid | social security privatization |
| democratic leader | committee on commerce science | **foreign oil** | billion trade deficit |
| federal spending | cord blood stem | president plan | asian pacific american |
| tax increase | medical liability reform | gun violence | president bush took office |

Other R: **saddam hussein, war on terrorism, iraqi people**

Other D: **funding for veterans health; war in iraq and afghanistan; improvised explosive device**

# Top Phrases: Fiscal Policy

| Republicans: 2-word | Republicans: 3-word | Democrats: 2-word | Democrats: 3-word |
|---|---|---|---|
| stem cell | embryonic stem cell | private accounts | veterans health care |
| natural gas | hate crimes legislation | trade agreement | congressional black caucus |
| **death tax** | adult stem cells | american people | va health care |
| illegal aliens | oil for food program | **tax breaks** | **billion in tax cuts** |
| class action | personal retirement accounts | trade deficit | credit card companies |
| war on terror | energy and natural resources | oil companies | security trust fund |
| embryonic stem | global war on terror | credit card | social security trust |
| **tax relief** | hate crimes law | nuclear option | privatize social security |
| illegal immigration | change hearts and minds | war in iraq | american free trade |
| date the time | global war on terrorism | middle class | central american free |
| boy scouts | class action fairness | african american | national wildlife refuge |
| hate crimes | committee on foreign relations | **budget cuts** | dependence on foreign oil |
| oil for food | **deficit reduction bill** | nuclear weapons | **tax cuts for the wealthy** |
| global war | boy scouts of america | checks and balances | vice president cheney |
| medical liability | **repeal of the death tax** | civil rights | arctic national wildlife |
| highway bill | highway trust fund | veterans health | bring our troops home |
| adult stem | action fairness act | **cut medicaid** | social security privatization |
| democratic leader | committee on commerce science | foreign oil | billion trade deficit |
| **federal spending** | cord blood stem | president plan | asian pacific american |
| **tax increase** | medical liability reform | gun violence | president bush took office |

Other R: **raise taxes; percent growth; increase taxes; growth rate; government spending; raising taxes; death tax repeal; million jobs created; percent growth rate**

Other D: **estate tax; budget deficit; bill cuts; medicaid cuts; cut funding; spending cuts; pay for tax cuts; cut student loans; cut food stamps; cut social security; billion in tax breaks**

# Obtain Phrase Counts from Newspapers

# Obtain Phrase Counts from Newspapers

ProQuest | **ProQuest Newsstand**

"personal retirement account" AND pub(washington times)

Full text                                                                    Modify search      Tips

**Suggested subjects**  Hide                                    Powered by ProQuest® Smart Search
Washington Times (Company/Org)      Washington Times (Company/Org) AND Newspapers
Washington Times (Company/Org) AND Moon, Sun Myung (Person)      Washington Times (Company/Org) AND Washington DC (Place)

**53 Results** *      Search within                              Create alert    Create RSS feed    Save search

Select 1-20   Brief view  |  Detailed view

**1**   Tin ears on Social Security: [2 Edition 1]
        Ferrara, Peter. **Washington Times** [Washington, D.C] 13 July 1999: A17.
        ...who wanted to advance a **personal retirement account** option to Social Security
        ...worker's wages into a **personal retirement account** for the worker. These payments
        ...and proposed a sound **personal retirement account** plan. He would have granted
        Citation/Abstract      Full text

**2**   Washington's financial miscreants sucker us
        Hurt, Charles. **Washington Times** [Washington, D.C] 05 Dec 2012: A.6.
        ...billions out of our **personal retirement account** to fund an obscene lavishness
        Citation/Abstract      Full text

**3**   Brickbats blur Bush proposal for Social Security ; Plan backers say 'truth' obscured
        Lambro, Donald. **Washington Times** [Washington, D.C] 06 Feb 2005: A03.
        ... Mr. Bush's **personal retirement account** (PRA) plan has been attacked by
        Citation/Abstract      Full text

**4**   Touching Social Security's hot third rail: [2 Edition]
        Lambro, Donald. **Washington Times** [Washington, D.C] 12 Dec 1996: A.17.
        ...Security taxes into their own **personal retirement account**. Mr. Forbes'
        Citation/Abstract      Full text

# Example: Social Security

- "House GOP offers plan for Social Security; Bush's **private accounts** would be scaled back" (*Washington Post*, 6/23/05)
- "GOP backs use of Social Security surplus ; Finds funding for **personal accounts**" (*Washington Times*, 6/23/05)

# Linear Model of Phrase Frequency

- Let $y_i$ be Republican vote share in senator $i$'s state
- Let $x_{ij}$ be share of senator $i's$ speech going to phrase $j$

$$E(x_{ij}|y_i) = \alpha_j + \beta_j y_i$$

- Estimate via least squares
  - Procedure called *marginal regression*
- Apply same model to newspapers to infer $y_m$

# Validation



Scatter plot with "Mondo Times conservativeness rating" on the x-axis (values 2, 3, 4) and "Slant index" on the y-axis (values .4, .45, .5).

Data points labeled:
- 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
- St. Petersburg Times
- Hackensack Record
- Orlando Sentinel
- Newark Star–Ledger
- Seattle Times
- Buffalo News
- Los Angeles Times
- Washington Post
- Boston Globe
- New York Times
- Philadelphia Inquirer
- USA Today
- Dallas Morning News
- Memphis Commercial Appeal
- Chicago Tribune
- San Francisco Chronicle
- New Orleans Times–Picayune
- Baltimore Sun
- Atlanta Constitution
- Detroit News
- Tri–Valley Herald

- Newspaper rankings consistent with external sources
- Phrases make sense
- Sensitivity analysis
  - Change scores $y_i$ (ADA, NOMINATE)
  - Change set of phrases
- Check agreement across sources
- Go look at the newspapers

TABLE A.I

AUDIT OF SEARCH RESULTS[a]

| | | | Share of Hits That Are | | | | | |
| Phrase | Total Hits | Share of Hits in Quotes | AP Wire Stories | Other Wire Stories | Letters to the Editor | Maybe Opinion | Clearly Opinion | Independently Produced News |
|---|---|---|---|---|---|---|---|---|
| Global war on terrorism | 2064 | 16% | 3% | 4% | 1% | 2% | 10% | 80% |
| Malpractice insurance | 2190 | 5% | 0% | 0% | 1% | 3% | 12% | 84% |
| Universal health care | 1523 | 9% | 1% | 0% | 7% | 8% | 28% | 56% |
| Assault weapons | 1411 | 9% | 3% | 12% | 4% | 1% | 25% | 56% |
| Child support enforcement | 1054 | 3% | 0% | 0% | 1% | 2% | 11% | 86% |
| Public broadcasting | 3375 | 8% | 1% | 0% | 2% | 4% | 22% | 71% |
| Death tax | 595 | 36% | 0% | 0% | 2% | 5% | 46% | 47% |
| Average (hit weighted) | | 10% | 1% | 2% | 3% | 3% | 19% | 71% |

[a] Authors' calculations based on ProQuest and NewsLibrary data base searches. See Appendix A for details.

# Economic Hypotheses

- Now that we have a measure, we can use it to model newspaper ideology
- Possible drivers
  - Consumer ideology
  - Owner ideology
  - Influence of incumbent politicians

# Role of Consumer Ideology

# Possible Confounds

- Reverse causality
- Slant is proxying for other newspaper attributes (e.g. emphasis on sports vs. business)
- Slant is proxying for other market attributes (e.g. geography)

# Soda vs. Pop



Generic names for Soft Drinks by county

# Solutions

- Control carefully for geography when relating slant to other variables
- Incorporate geography into predictive model
  - Predict the component of congressperson ideology that is orthogonal to Census division

- Can use predictive modeling as an aid to social science
- But
  - You get out what you put in
  - Consider possible sources of bias and misspecification

# Taddy (2013)

- Two main limitations of Gentzkow & Shapiro (2010)
  - Feature selection separate from model estimation
  - Linear model doesn't exploit multinomial structure of data

# Taddy (2013)

- Let $y_i$ be Republican vote share in senator $i$'s state
- Let $x_{ijt}$ be an indicator for senator $i$ says phrase $j$ at occasion $t$
- Then
$$\Pr\left(x_{ij} = 1\right) = \frac{\exp\left(\alpha_j + \beta_j y_i\right)}{\sum_{j'} \exp\left(\alpha_{j'} + \beta_{j'} y_i\right)}$$
- Estimate via maximum likelihood
  - Uses log penalty for regularization
  - Uses novel algorithm for maximization
- Penalty imposes sparsity in $\beta_j$s
  - Means we can use a very large number of phrases $p$
  - Can think of this as a way to maximize performance subject to a phrase "budget"

# Taddy (2013)



**Political Speech**

- Note: LDA = latent Dirichlet allocation (stay tuned)

109th Congress Vote–Shares

# Sentiment Analysis: *Financial News*

- Questions
  - What explains time-series/cross-section of equity returns?
  - Is there information beyond what is reflected in quantitative fundamentals (e.g. earnings)?

- Data: counts of words in WSJ "Abreast of the Market" column

- Features: Counts of words in each of 77 "Harvard-IV General Inquirer" categories
  - Weak
  - Positive
  - Negative
  - Active
  - Passive
  - etc.

## List of entries in tag category:

**Weak**

**List shows first 100 entries. Total number of entries in this category:**

**755**

**Entries for this category are shown with all tags assigned and sense definitions:**

**ABANDON**

H4Lvd Negativ Ngtv Weak Fail IAV AffLoss AffTot SUPV

**ABANDONMENT**

H4 Negativ Weak Fail Noun

**ABDICATE**

H4 Negativ Weak Submit Passive Finish IAV SUPV

- Regressors
  - Weak words
  - Negative words
  - First principal component ("pessimism")

**Table II**

**Predicting Dow Jones Returns Using Negative Sentiment**

The table data come from CRSP, NYSE, and the General Inquirer program. This table shows OLS estimates of the coefficient $\gamma_1$ in equation (1). Each coefficient measures the impact of a one-standard deviation increase in negative investor sentiment on returns in basis points (one basis point equals a daily return of 0.01%). The regression is based on 3,709 observations from January 1, 1984, to September 17, 1999. I use Newey and West (1987) standard errors that are robust to heteroskedasticity and autocorrelation up to five lags. Bold denotes significance at the 5% level; italics and bold denotes significance at the 1% level.

| News Measure | Regressand: Dow Jones Returns | | |
| --- | --- | --- | --- |
| | Pessimism | Negative | Weak |
| $BdNws_{t-1}$ | *−8.1* | *−4.4* | *−6.0* |
| $BdNws_{t-2}$ | 0.4 | 3.6 | 2.0 |
| $BdNws_{t-3}$ | 0.5 | −2.4 | −1.2 |
| $BdNws_{t-4}$ | **4.7** | **4.4** | *6.3* |
| $BdNws_{t-5}$ | 1.2 | 2.9 | **3.6** |
| $\chi^2(5)$ [*Joint*] | *20.0* | *20.8* | *26.5* |
| $p$-value | 0.001 | 0.001 | 0.000 |
| Sum of 2 to 5 | **6.8** | *9.5* | *10.7* |
| $\chi^2(1)$ [*Reversal*] | **4.05** | *8.35* | *10.1* |
| $p$-value | 0.044 | 0.004 | 0.002 |

# Antweiler and Frank (2004)

- Data: Message board contents on Yahoo! Finance and Raging Bull

```
--------------------
FROM YF
COMP ETYS
MGID 13639
NAME CaptainLihai
LINK 1
DATE 2000/01/25 04:11
SKIP
TITL ETYS will surprise all pt II
SKIP
TEXT ETYS will surprise all when it drops to below 15$ a pop, and even then
TEXT it will be too expensive.
TEXT
TEXT If the DOJ report is real, there will definately be a backlash against
TEXT the stock. Watch your asses. Get out while you can.
--------------------
FROM YF
COMP IBM
MGID 43653
NAME plainfielder
LINK 1
DATE 2000/03/29 11:39
SKIP
TITL BUY ON DIPS - This is the opportunity
SKIP
TEXT to make $$$ when IBM will be going up again following this profit taking
TEXT bout by Abbey Cohen and her brokerage firm.
TEXT
TEXT IBM shall go up again after today.
-----------------
```

# Antweiler and Frank (2004)

- Count words
- Create training set of 1000 messages hand-coded as buy, sell, hold
- Compute "naive Bayes classification:" posterior guess assuming words are independent

**Table I**
**Naive Bayes Classification Accuracy within Sample and Overall Classification Distribution**

The first percentage column shows the actual shares of 1,000 hand-coded messages that were classified as buy (B), hold (H), or sell (S). The buy-hold-sell matrix entries show the in-sample prediction accuracy of the classification algorithm with respect to the learned samples, which were classified by the authors (Us).

| Classified: by Us | % | By Algorithm | | |
|---|---|---|---|---|
| | | Buy | Hold | Sell |
| Buy | 25.2 | 18.1 | 7.1 | 0.0 |
| Hold | 69.3 | 3.4 | 65.9 | 0.0 |
| Sell | 5.5 | 0.2 | 1.2 | 4.1 |
| 1,000 messages[a] | | 21.7 | 74.2 | 4.1 |
| All messages[b] | | 20.0 | 78.8 | 1.3 |

[a]These are the 1,000 messages contained in the training data set.
[b]This line provides summary statistics for the out-of-sample classification of all 1,559,621 messages.

# Antweiler and Frank (2004)

- Small amount of predictability in returns
- Messages predict volatility
- Disagreement (variable recommendations) predicts volume

# Other Examples

- Li (2010): Uses naive Bayes to measure sentiment of forward-looking statements in 10Ks/10Qs
- Hanley and Hoberg (2012): Use cosine distance to measure revisions to IPO prospectuses

# Topic Models

# Factor Models

- "Unsupervised" methods (factor analysis, PCA) project high-dimensional data into low-dimensional measures, preserving as much variation as possible.

- "Unsupervised" methods (factor analysis, PCA) project high-dimensional data into low-dimensional measures, preserving as much variation as possible.

- E.g.,
    - Congressional roll call votes $\rightarrow$"Common space" scores
    - Survey responses $\rightarrow$ "Big 5" personality traits

- "Unsupervised" methods (factor analysis, PCA) project high-dimensional data into low-dimensional measures, preserving as much variation as possible.

- E.g.,
  - Congressional roll call votes $\rightarrow$ "Common space" scores
  - Survey responses $\rightarrow$ "Big 5" personality traits

- Low dimensional measures are then inputs into subsequent analysis

# Factor Models

- "Unsupervised" methods (factor analysis, PCA) project high-dimensional data into low-dimensional measures, preserving as much variation as possible.
- E.g.,
  - Congressional roll call votes $\rightarrow$ "Common space" scores
  - Survey responses $\rightarrow$ "Big 5" personality traits
- Low dimensional measures are then inputs into subsequent analysis
- E.g.,
  - How has polarization in Congress changed over time? (Poole & Rosenthal 1984)
  - How does personality correlate with job performance (Tett et al. 1991)

- Topic models extend these methods to multinomial data such as text
- Relevant to measuring, e.g.,
  - What people talk about on social networks
  - What products share similar descriptions on Amazon / EBay
  - What "stories" are in the news today
  - What are economists studying

- As with other unsupervised methods, topic models are of most interest to social scientists as an input into subsequent analysis

# Purpose

- As with other unsupervised methods, topic models are of most interest to social scientists as an input into subsequent analysis
- E.g.,
  - Do discussions of particular topics on Twitter predict stock movements?
  - Which products are close substitutes on EBay?
  - Is media slant driven by what you talk about or how you talk about it?
  - How has the distribution of topics in economics changed over time?

- As with other unsupervised methods, topic models are of most interest to social scientists as an input into subsequent analysis
- E.g.,
  - Do discussions of particular topics on Twitter predict stock movements?
  - Which products are close substitutes on EBay?
  - Is media slant driven by what you talk about or how you talk about it?
  - How has the distribution of topics in economics changed over time?
- A fair critique of topic modeling literature is that it hasn't progressed much beyond the measurement stage

# Topic Models: *Blei & Lafferty (2006)*

- OCR text of *Science* 1880-2002 (from JSTOR)
  - Count words used 25 or more times (after stemming and removing stopwords)
  - Vocabulary: $15,955$ words
  - Total documents: $30,000$ articles

# Output

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| human | evolution | disease | computer |
| genome | evolutionary | host | models |
| dna | species | bacteria | information |
| genetic | organisms | diseases | data |
| genes | life | resistance | computers |
| sequence | origin | bacterial | system |
| gene | biology | new | network |
| molecular | groups | strains | systems |
| sequencing | phylogenetic | control | model |
| map | living | infectious | parallel |
| information | diversity | malaria | methods |
| genetics | group | parasite | networks |
| mapping | new | parasites | software |
| project | two | united | new |
| sequences | common | tuberculosis | simulations |

# Output



**"Theoretical Physics"**

**"Neuroscience"**

- Latent Dirichlet Allocation (LDA) introduced by Blei et al. (2003) as an extension of factor models to discrete data

- Setup
  - Documents $i \in \{1, ..., n\}$
  - Words $j \in \{1, ..., p\}$
  - Data $\mathbf{x}_i$ is $(1 \times p)$ vector of word counts for document $i$

# Model: LDA

- Setup
  - Documents $i \in \{1, ..., n\}$
  - Words $j \in \{1, ..., p\}$
  - Data $\mathbf{x}_i$ is $(1 \times p)$ vector of word counts for document $i$

- Factor model
  - $\theta_{ik}$ is value of $k$-th **factor** for document $i$
  - $\boldsymbol{\beta}_k$ is $(1 \times p)$ vector of **loadings** for factor $k$

$$E(\mathbf{x}_i) = \boldsymbol{\beta}_1 \theta_{i1} + ... \boldsymbol{\beta}_K \theta_{iK}$$

# Model: LDA

- Setup
  - Documents $i \in \{1, ..., n\}$
  - Words $j \in \{1, ..., p\}$
  - Data $\mathbf{x}_i$ is $(1 \times p)$ vector of word counts for document $i$

- Factor model
  - $\theta_{ik}$ is value of $k$-th **factor** for document $i$
  - $\boldsymbol{\beta_k}$ is $(1 \times p)$ vector of **loadings** for factor $k$

  $$E\left(\mathbf{x}_i\right) = \boldsymbol{\beta_1}\theta_{i1} + ...\boldsymbol{\beta_K}\theta_{iK}$$

- LDA
  - $\theta_{ik}$ is weight on $k$-th **topic** for document $i$
  - $\boldsymbol{\beta_k}$ is $(1 \times p)$ vector of **word probabilities** for topic $k$

  $$\mathbf{x}_i \sim Multinomial\left(\boldsymbol{\beta_1}\theta_{i1} + ...\boldsymbol{\beta_K}\theta_{iK}\right)$$

**Seeking Life's Bare (Genetic) Necessities**

COLD SPRING HARBOR, NEW YORK— How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

# Model: LDA

# Model: LDA

- For each document $i$...
  - Draw topic proportions $\theta_i$ from Dirichlet distribution with parameter $\alpha$
  - For each word $j$...
    - Draw a topic assignment $k \sim Multinomial(\theta_i)$
    - Draw word $x_{ij} \sim Multinomial(\beta_k)$

# Model: Dynamic

- One limitation of LDA is it assumes documents are exchangeable; in many settings of interest, topics evolve systematically over time

# Model: Dynamic

## "Instantaneous Photography" (1890)



## "Infrared Reflectance in Leaf-Sitting Neotropical Frogs" (1977)

- Divide text into sequential slices (e.g., by year)
- Assume each slice's documents drawn from LDA model
- Allow word distribution within topics $\beta$ and distribution over topics $\alpha$ to evolve via markov process

# Estimation

- Bayesian inference intractable using standard methods (e.g., Gibbs sampling)
  - Blei (2006) $\rightarrow$ variational inference
  - Taddy R package $\rightarrow$ MAP estimation
  - Current favorite $\rightarrow$ Stochastic gradient descent
- Main estimates are for 20 topic model

# Results: LDA

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| human | evolution | disease | computer |
| genome | evolutionary | host | models |
| dna | species | bacteria | information |
| genetic | organisms | diseases | data |
| genes | life | resistance | computers |
| sequence | origin | bacterial | system |
| gene | biology | new | network |
| molecular | groups | strains | systems |
| sequencing | phylogenetic | control | model |
| map | living | infectious | parallel |
| information | diversity | malaria | methods |
| genetics | group | parasite | networks |
| mapping | new | parasites | software |
| project | two | united | new |
| sequences | common | tuberculosis | simulations |

# Results: Dynamic

| **1880** | **1890** | **1900** | **1910** | **1920** | **1930** | **1940** |
|---|---|---|---|---|---|---|
| electric | electric | apparatus | air | apparatus | tube | air |
| machine | power | steam | water | tube | apparatus | tube |
| power | company | power | engineering | air | glass | apparatus |
| engine | steam | engine | apparatus | pressure | air | glass |
| steam | electrical | engineering | room | water | mercury | laboratory |
| two | machine | water | laboratory | glass | laboratory | rubber |
| machines | two | construction | engineer | gas | pressure | pressure |
| iron | system | engineer | made | made | made | small |
| battery | motor | room | gas | laboratory | gas | mercury |
| wire | engine | feet | tube | mercury | small | gas |

| **1950** | **1960** | **1970** | **1980** | **1990** | **2000** |
|---|---|---|---|---|---|
| tube | tube | air | high | materials | devices |
| apparatus | system | heat | power | high | device |
| glass | temperature | power | design | power | materials |
| air | air | system | heat | current | current |
| chamber | heat | temperature | system | applications | gate |
| instrument | chamber | chamber | systems | technology | high |
| small | power | high | devices | devices | light |
| laboratory | high | flow | instruments | design | silicon |
| pressure | instrument | tube | control | device | material |
| rubber | control | design | large | heat | technology |

# Results: Dynamic

The Brain of the Orang (1880)

## Representation of the Visual Field on the Medial Wall of Occipital-Parietal Cortex in the Owl Monkey (1976)

# Topic Models: *Quinn et al. (2010)*

- Full text of speeches in US Senate 1995-2004
  - Count words appearing in 0.5% or more of speeches (after stemming)
  - Vocabulary: $3,807$ words
  - Total documents: $118,065$ speeches

- Like Blei & Lafferty (2006), except
  1. Each document is in exactly one topic
  2. Dynamic distribution of topics, but topics themselves are static

# Model

- Blei & Lafferty (2006)
  - $\mathbf{x}_i \sim Multinomial\left(\boldsymbol{\beta}_1 \theta_{i1} + ... \boldsymbol{\beta}_K \theta_{iK}\right)$
  - $\boldsymbol{\theta}_i \sim F\left(\alpha\right)$
  - $\beta$ and $\alpha$ both evolve over time

- Quinn et al. (2010)
  - $\mathbf{x}_i \sim Multinomial\left(\beta_{k(i)}\right)$
  - $Pr\left(k\left(i\right) = j\right) = \alpha_j$
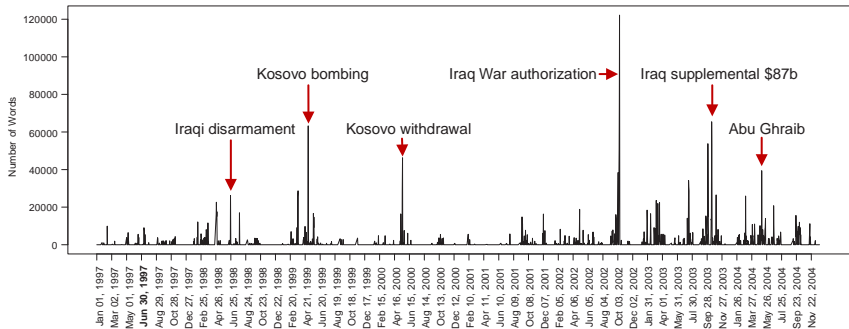  - $\alpha$ evolves over time; $\beta$ constant

# Estimation

- Estimate using ECM algorithm
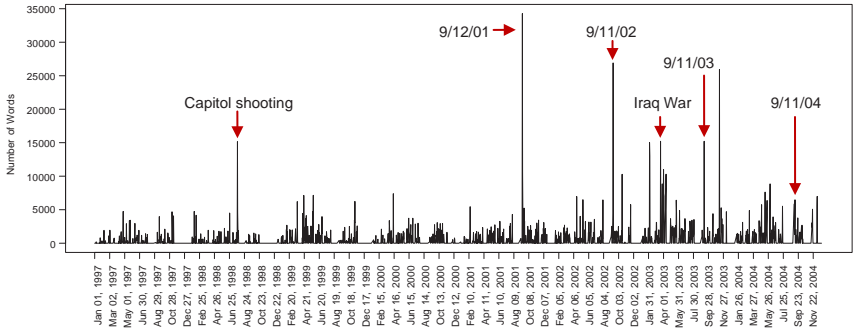- Main estimates are for 42 topic model (chosen based on "substantive and conceptual" criteria)
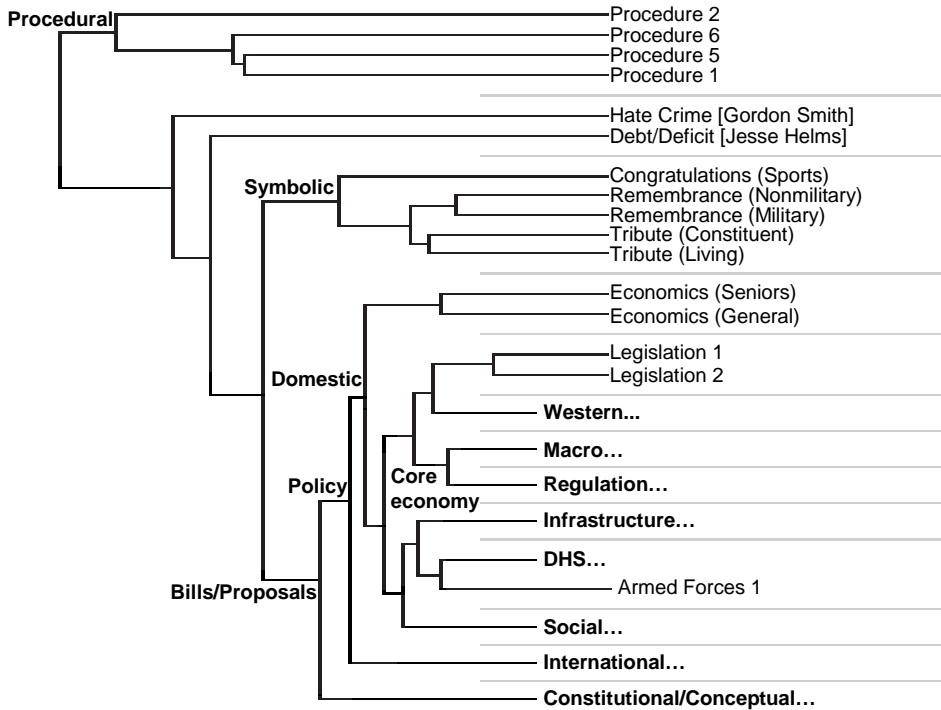
**TABLE 3  Topic Keywords for 42-Topic Model**

| Topic (Short Label) | Keys |
|---|---|
| 1. Judicial Nominations | *nomine, confirm, nomin, circuit, hear, court, judg, judici, case, vacanc* |
| 2. Constitutional | *case, court, attorney, supreme, justic, nomin, judg, m, decis, constitut* |
| 3. Campaign Finance | *campaign, candid, elect, monei, contribut, polit, soft, ad, parti, limit* |
| 4. Abortion | *procedur, abort, babi, thi, life, doctor, human, ban, decis, or* |
| 5. Crime 1 [Violent] | *enforc, act, crime, gun, law, victim, violenc, abus, prevent, juvenil* |
| 6. Child Protection | *gun, tobacco, smoke, kid, show, firearm, crime, kill, law, school* |
| 7. Health 1 [Medical] | *diseas, cancer, research, health, prevent, patient, treatment, devic, food* |
| 8. Social Welfare | *care, health, act, home, hospit, support, children, educ, student, nurs* |
| 9. Education | *school, teacher, educ, student, children, test, local, learn, district, class* |
| 10. Military 1 [Manpower] | *veteran, va, forc, militari, care, reserv, serv, men, guard, member* |
| 11. Military 2 [Infrastructure] | *appropri, defens, forc, report, request, confer, guard, depart, fund, project* |
| 12. Intelligence | *intellig, homeland, commiss, depart, agenc, director, secur, base, defens* |
| 13. Crime 2 [Federal] | *act, inform, enforc, record, law, court, section, crimin, internet, investig* |
| 14. Environment 1 [Public Lands] | *land, water, park, act, river, natur, wildlif, area, conserv, forest* |
| 15. Commercial Infrastructure | *small, busi, act, highwai, transport, internet, loan, credit, local, capit* |
| 16. Banking / Finance | *bankruptci, bank, credit, case, ir, compani, file, card, financi, lawyer* |
| 17. Labor 1 [Workers] | *worker, social, retir, benefit, plan, act, employ, pension, small, employe* |

**Defense [Use of Force]**

**Symbolic [Remembrance − Military]**

# Conclusion

# Conclusion

- Today
  - Prediction with high-dimensional data
  - Applications

- Tomorrow
  - Estimating treatment effects with high-dimensional data
  - Application