

# SAMUEL SHERMAN



You can find me at:

[github.com/scsherm](https://github.com/scsherm)

[linkedin.com/in/samuelcsherman](https://linkedin.com/in/samuelcsherman)

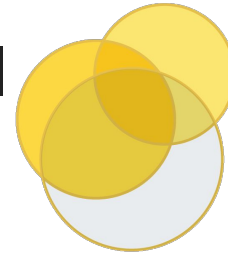
# Congressional Bill Modeling

# Data Collection

- ◎ **Congress.gov**
  - Bill text

CONGRESS.GOV

- ◎ **Sunlight Foundation API**
  - Bill data (JSON)
  - Votes data (JSON)



SUNLIGHT  
FOUNDATION

- ◎ **MongoDB**



mongoDB

- ◎ **Tools**
  - Pandas, pymongo, sklearn, nltk, numpy, scipy, AWS ec2

## Motivations

- ◉ **Percent of yes votes for a given party**
  - **Polarization in government**
- ◉ **Whether a bill will reach a vote**
  - **Important features in bills**
- ◉ **Latent topics in bill text**
  - **Prevalence over time**
  - **Distinguishable characteristics between presidencies**

# Natural Language Processing

## ◎ Stopwords

- Regular words: “and”, “the”, “be”, “it”, “there”
- Congressional: “amendment”, “bill”, “quorum”, “act”

## ◎ Maximum document frequency

## ◎ TFIDF

- TF: Frequency of words across a single document
- IDF: Frequency of words across all documents

# Non-Negative Matrix Factorization

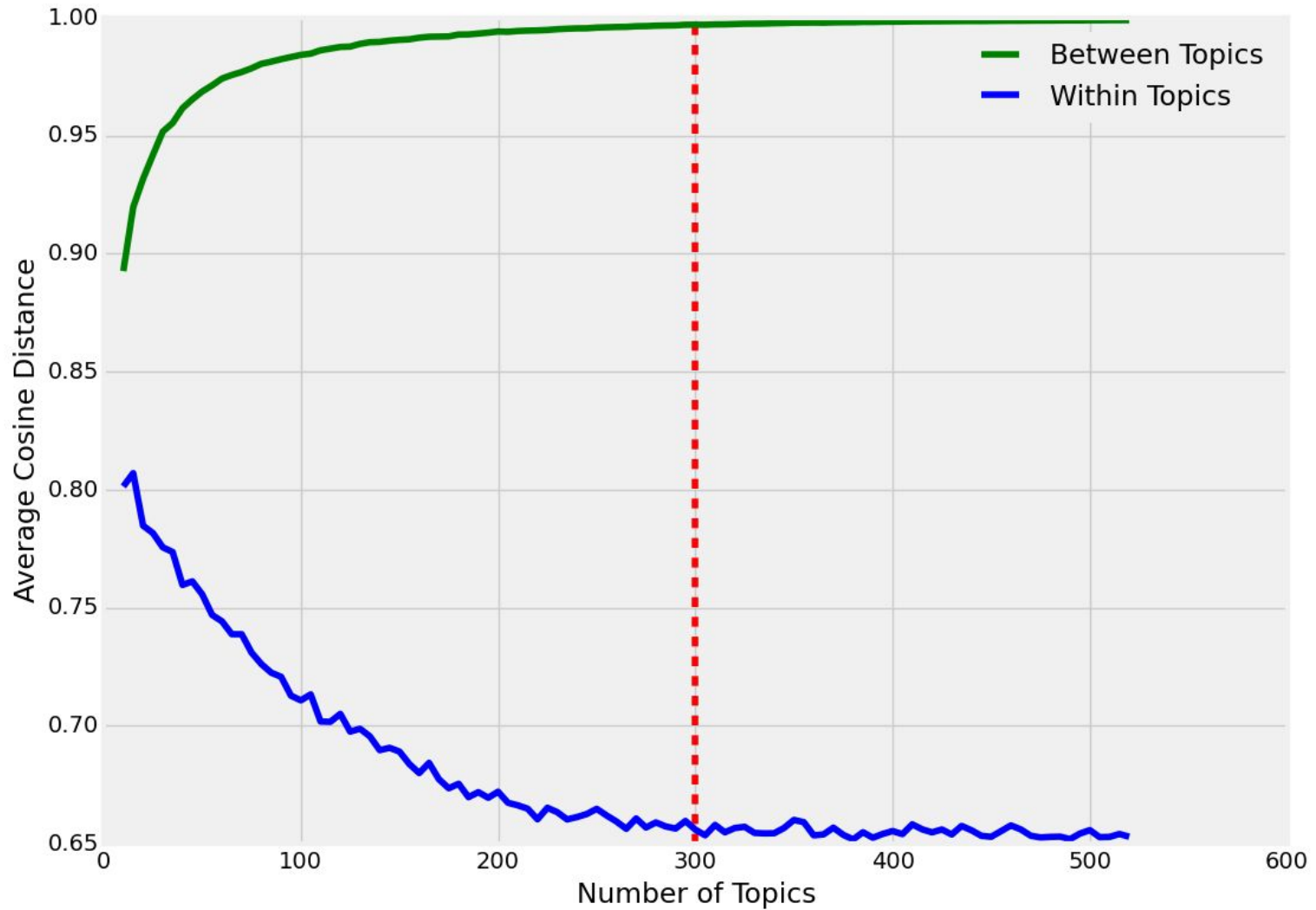
- **Similarity within topics**  **high**



- **Similarity between topics**  **low**

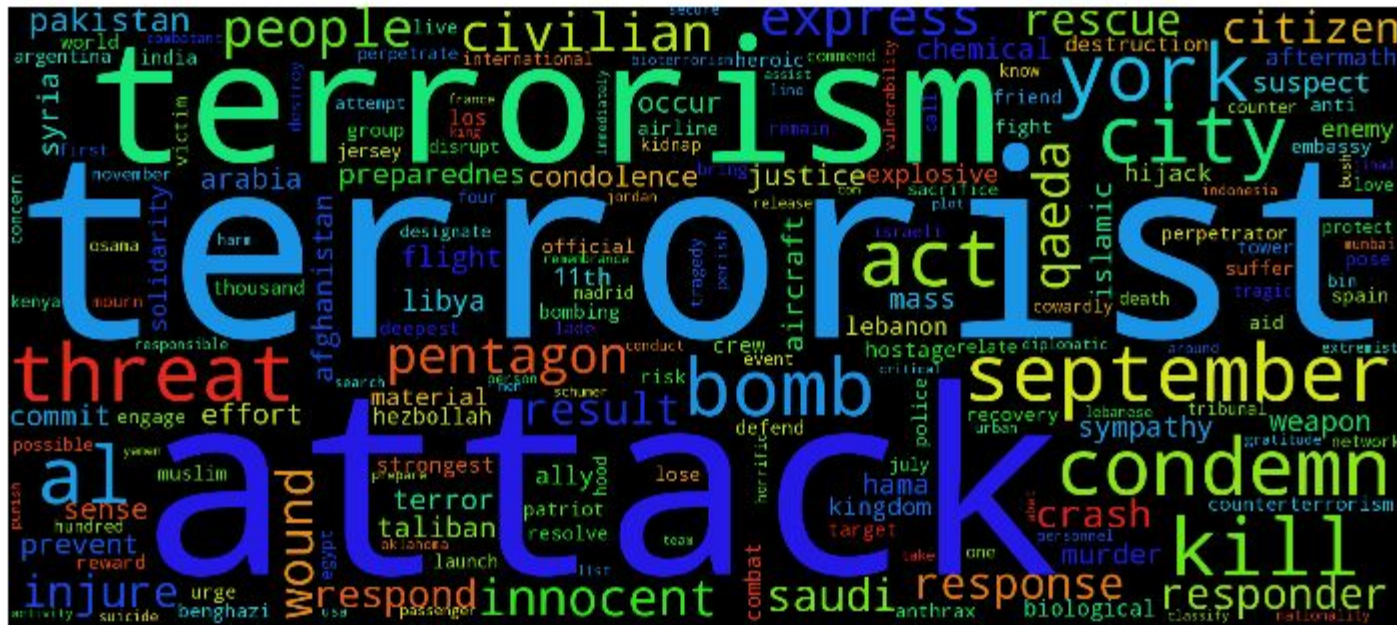


# Choosing K Topics



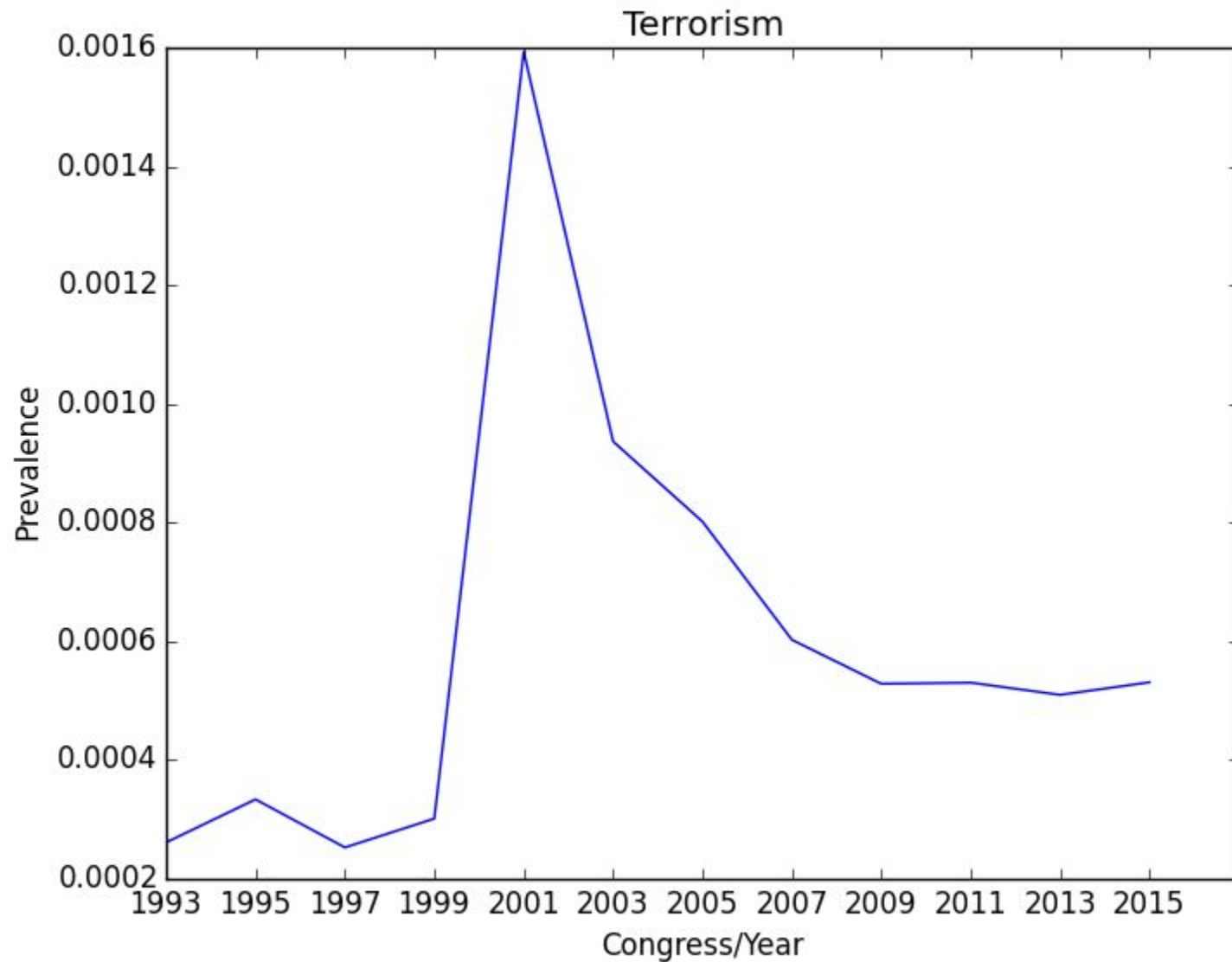
*Measuring Topic Distance*

# Topic 110



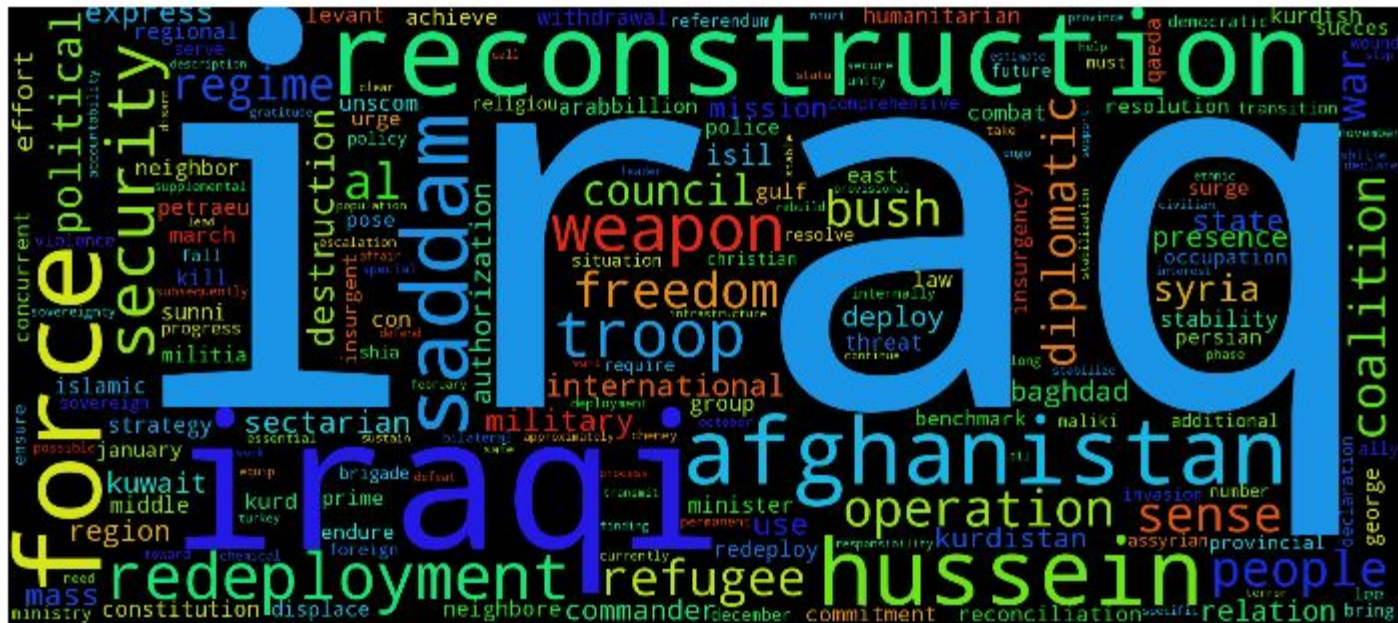
# Terrorism



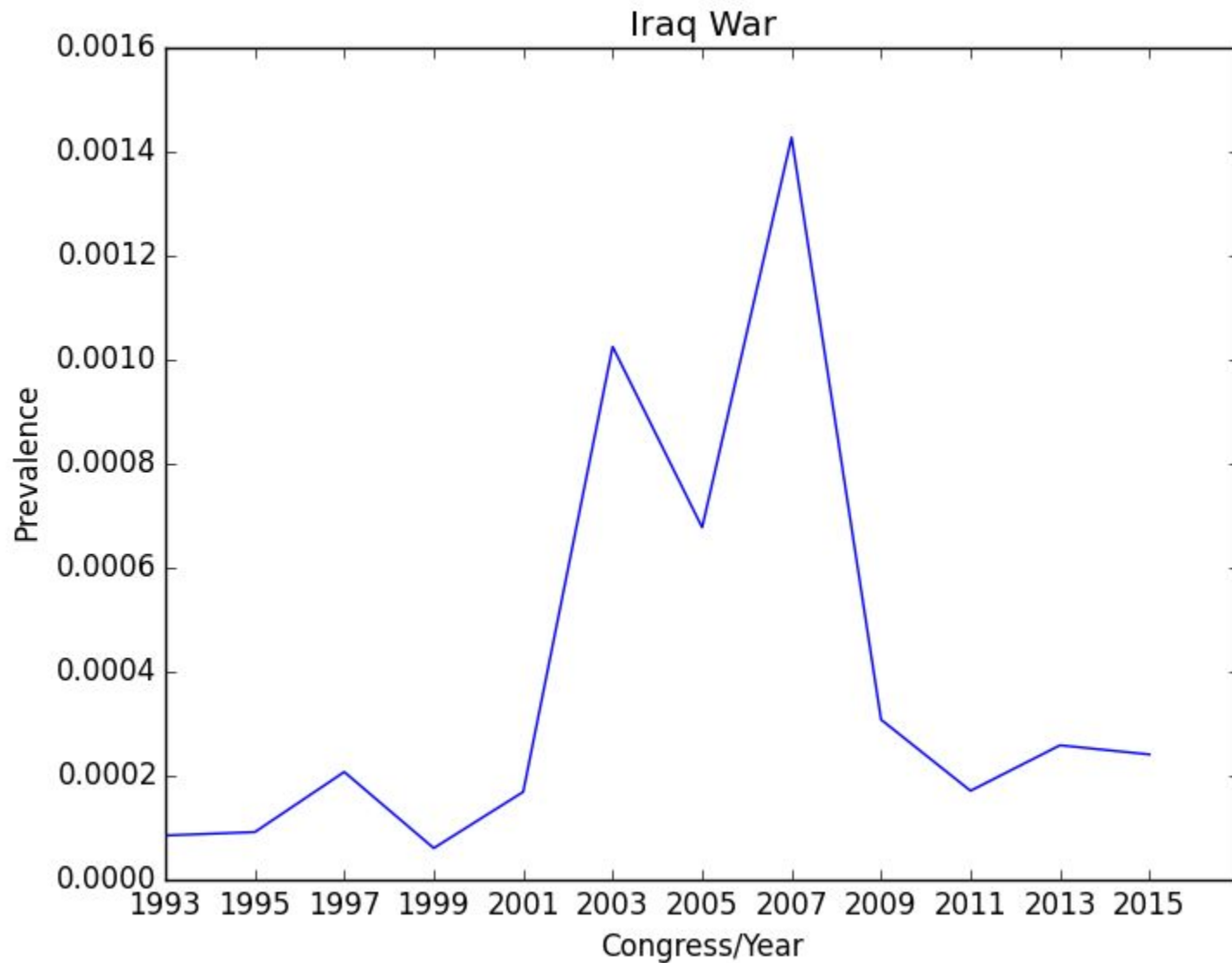


### *Terrorism in Bills Over Time*

# Topic 76

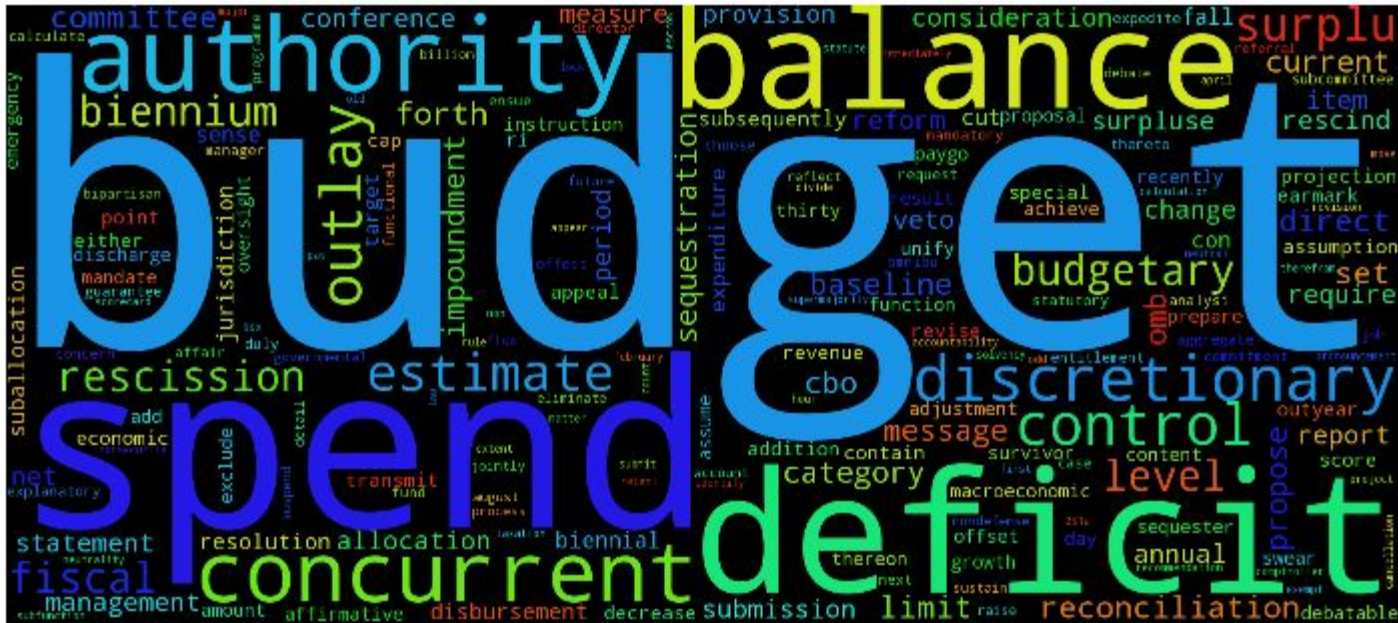


## *Iraq War*

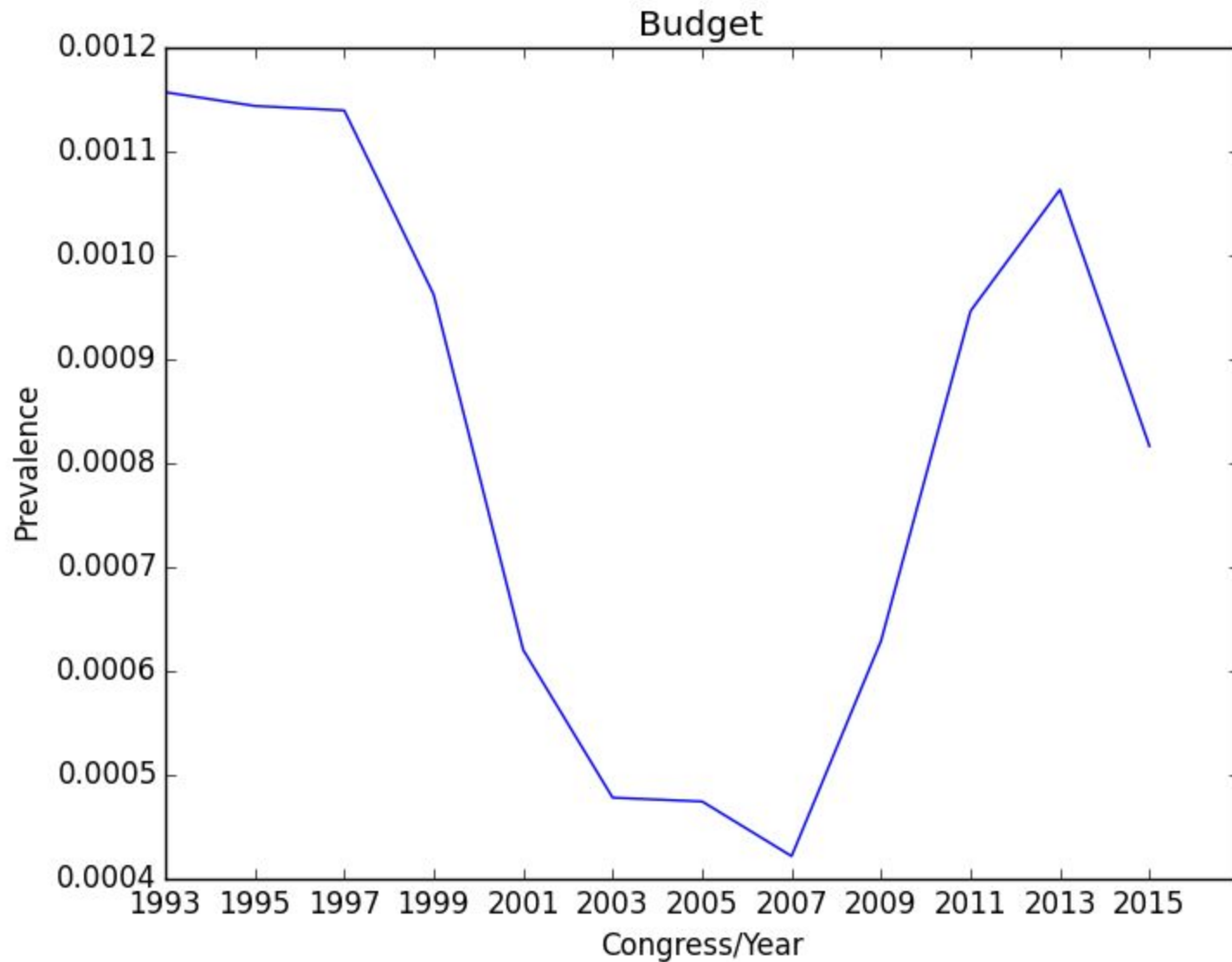


### *Iraq War in Bills Over Time*

# Topic 188



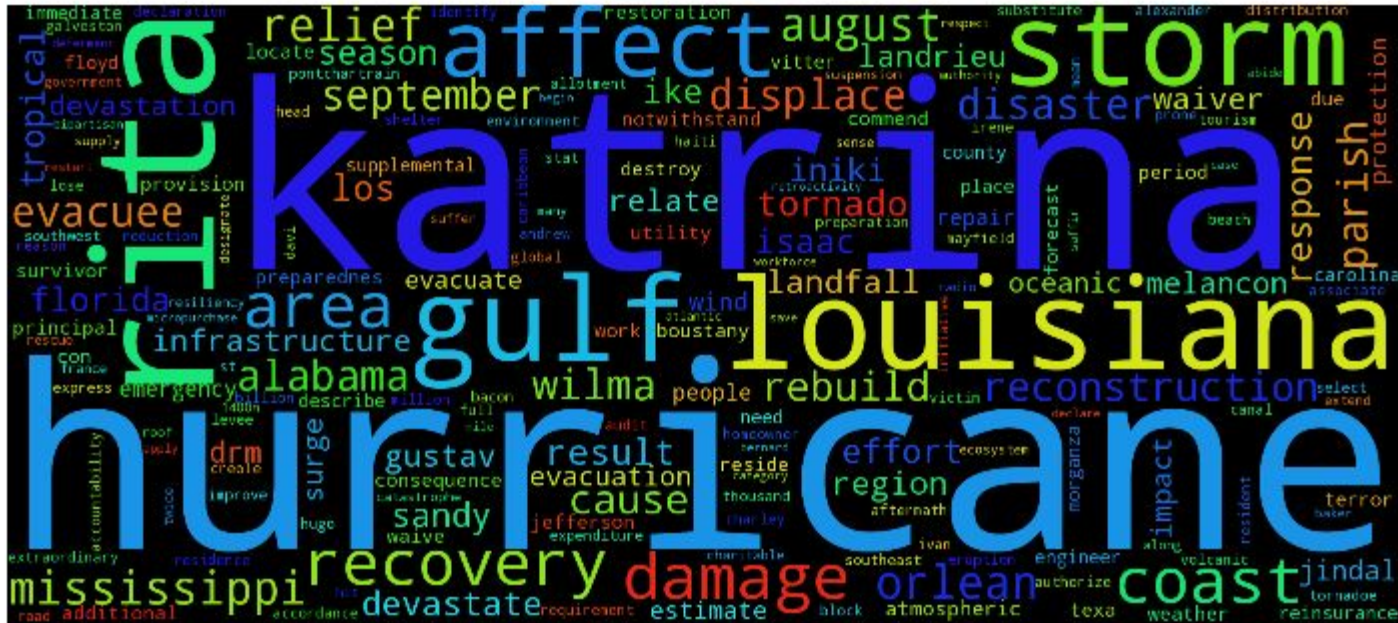
## Budget



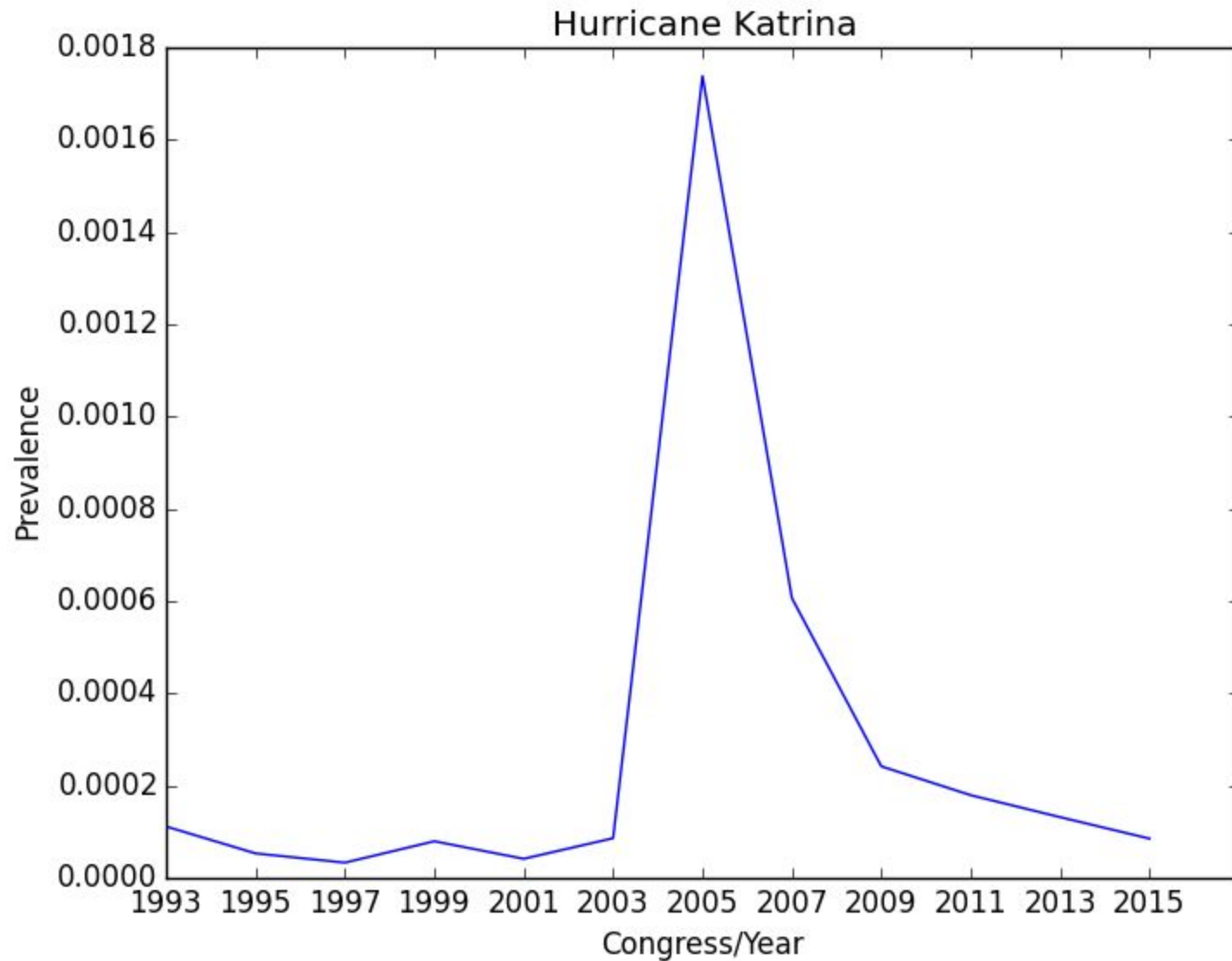
### *Budget in Bills Over Time*



# Topic 205

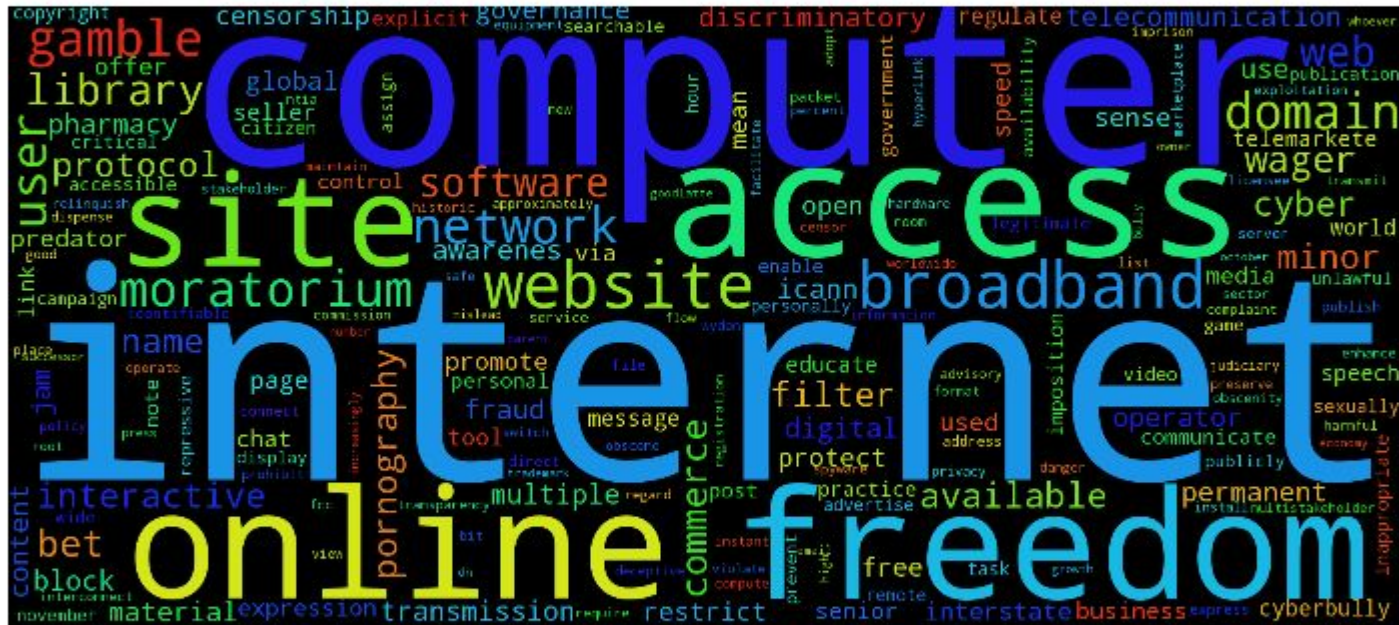


# Hurricane Katrina



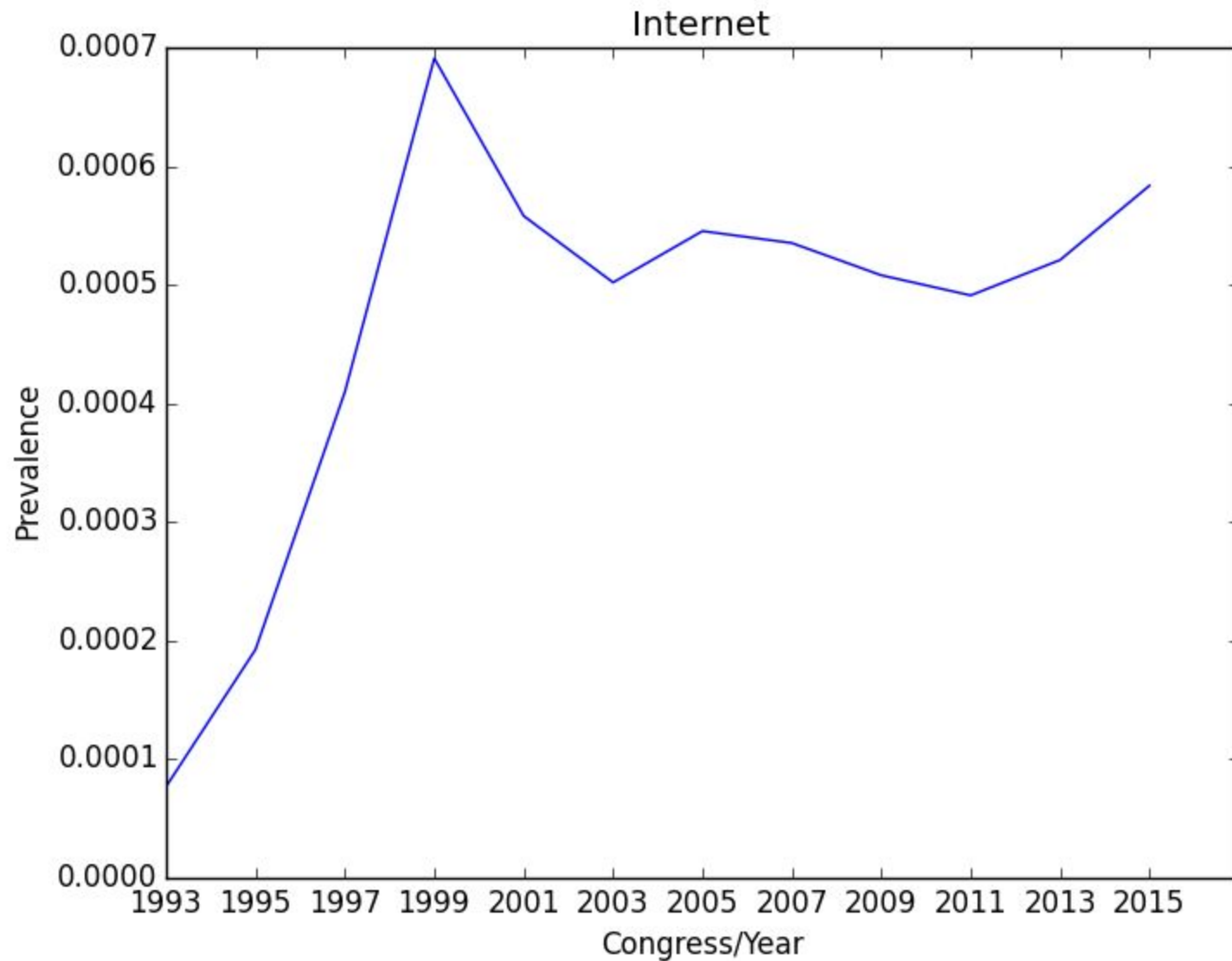
### *Hurricane Katrina in Bills Over Time*

# Topic 220



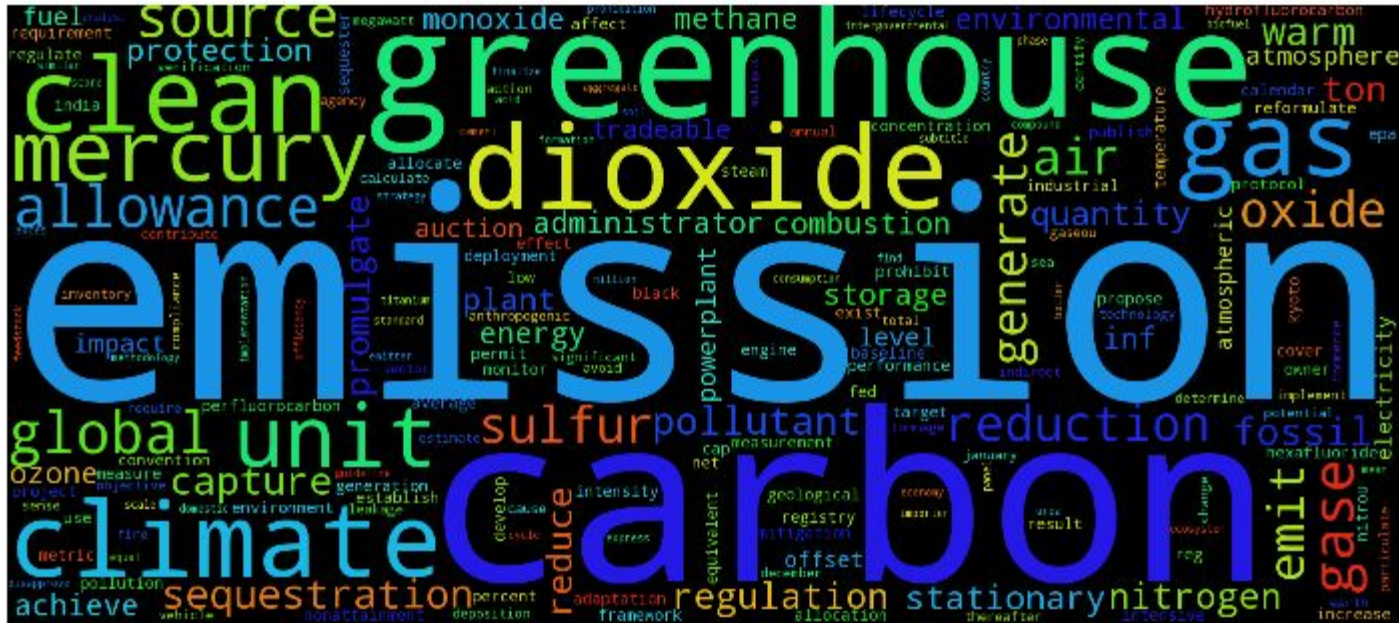
## Internet



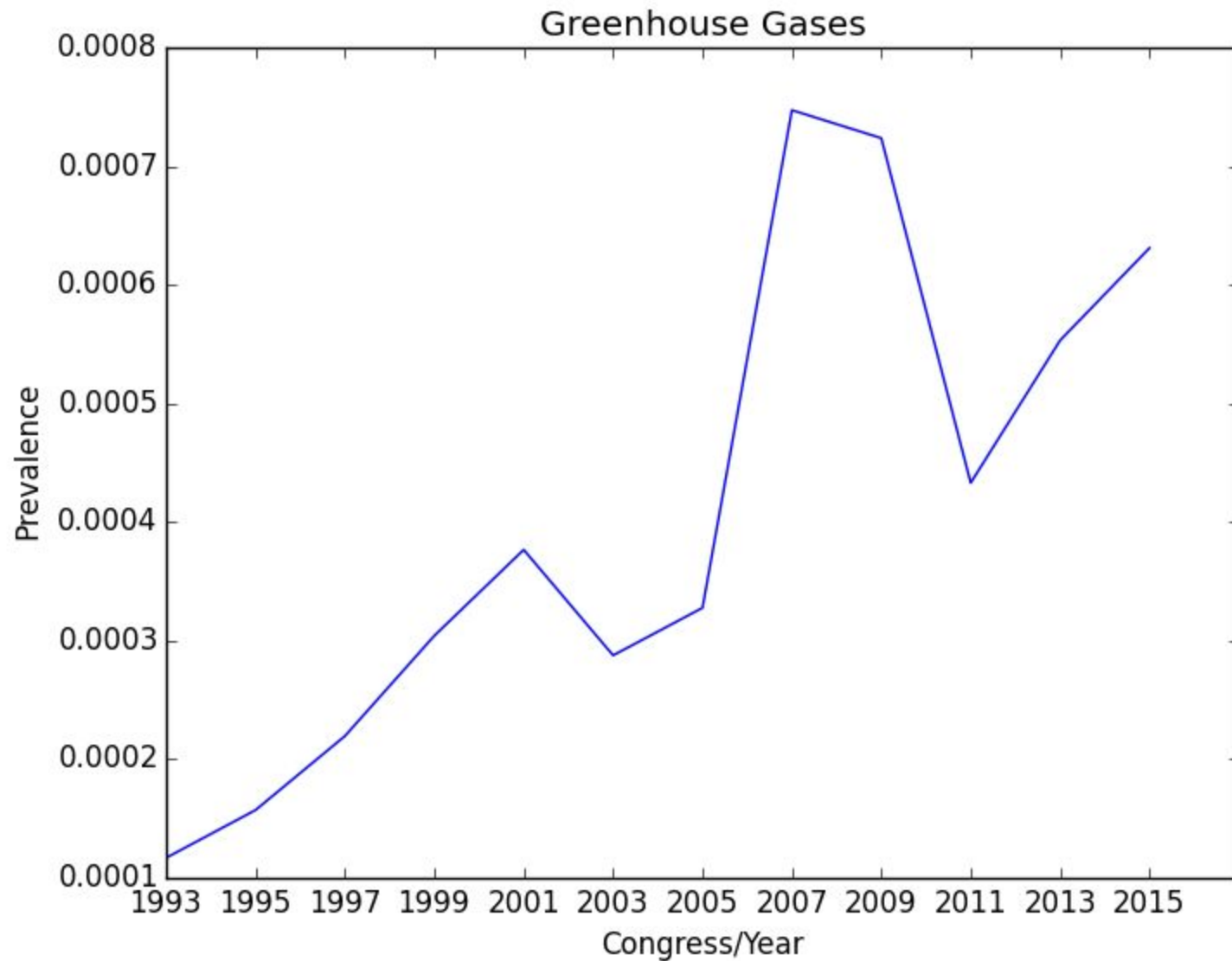


### *Internet in Bills Over Time*

# Topic 277

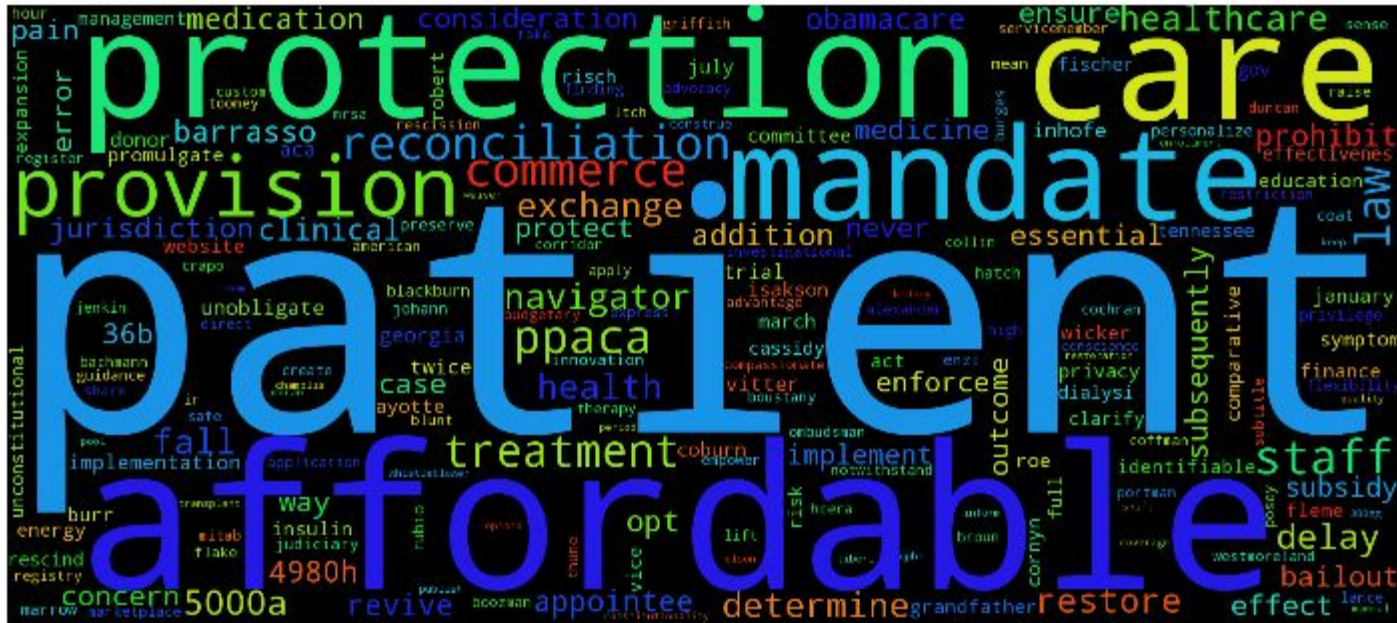


# Greenhouse Gases

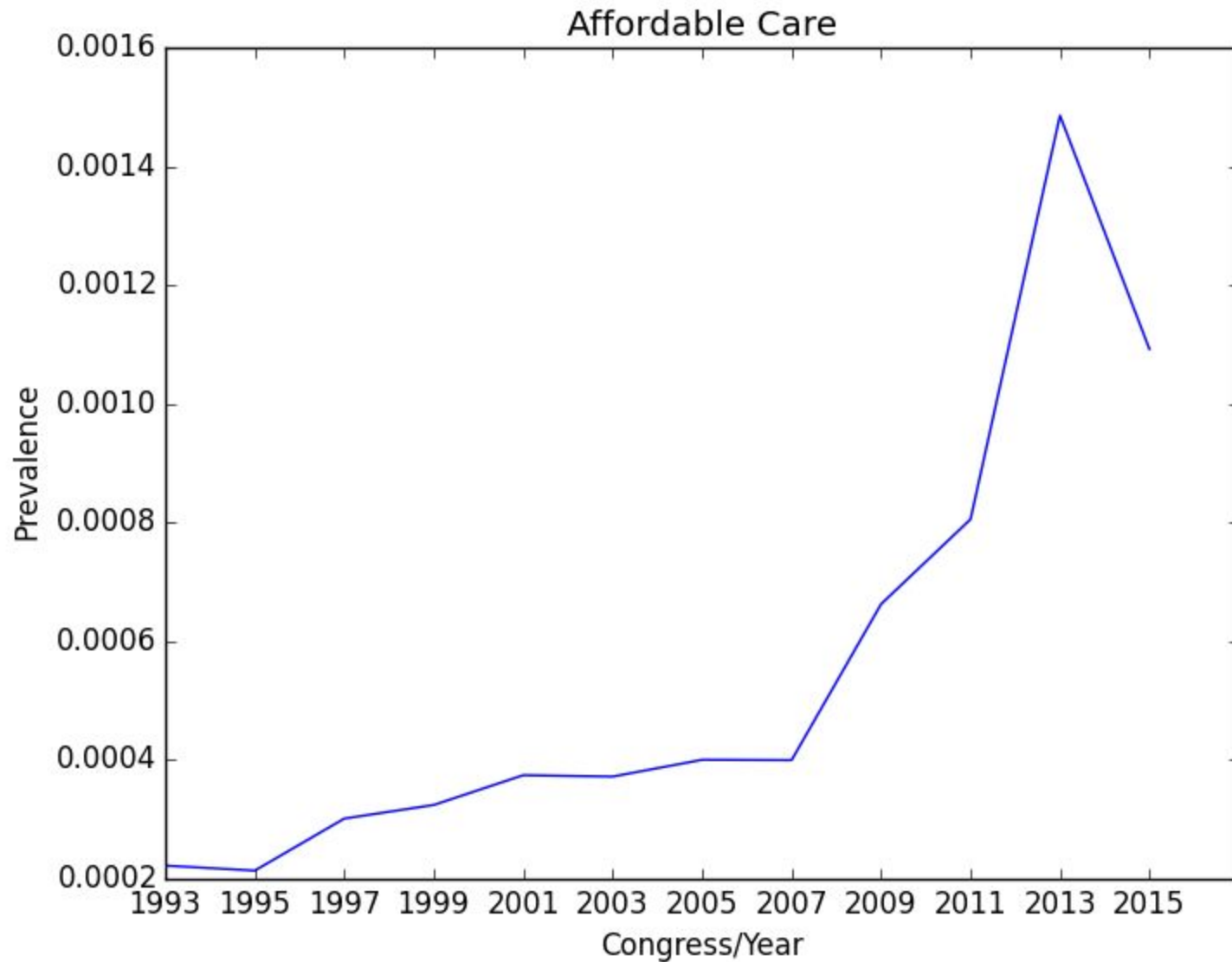


## *Greenhouse Gases in Bills Over Time*

# Topic 286



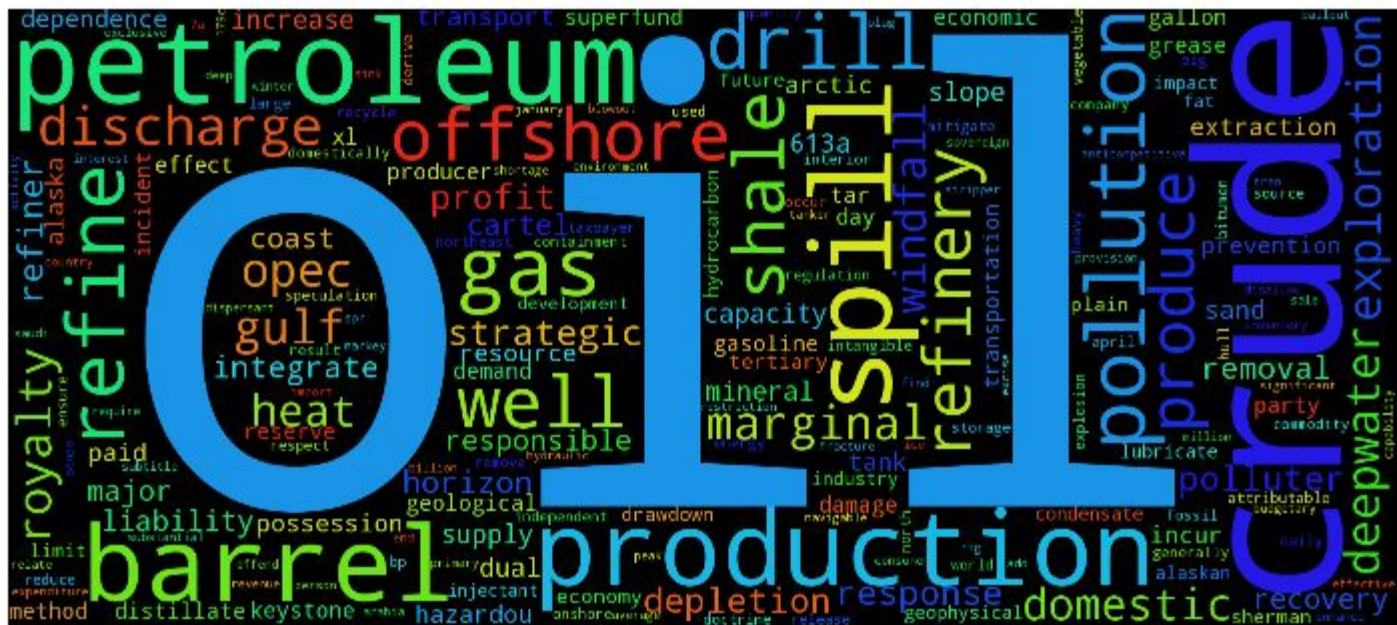
## Affordable Care Act



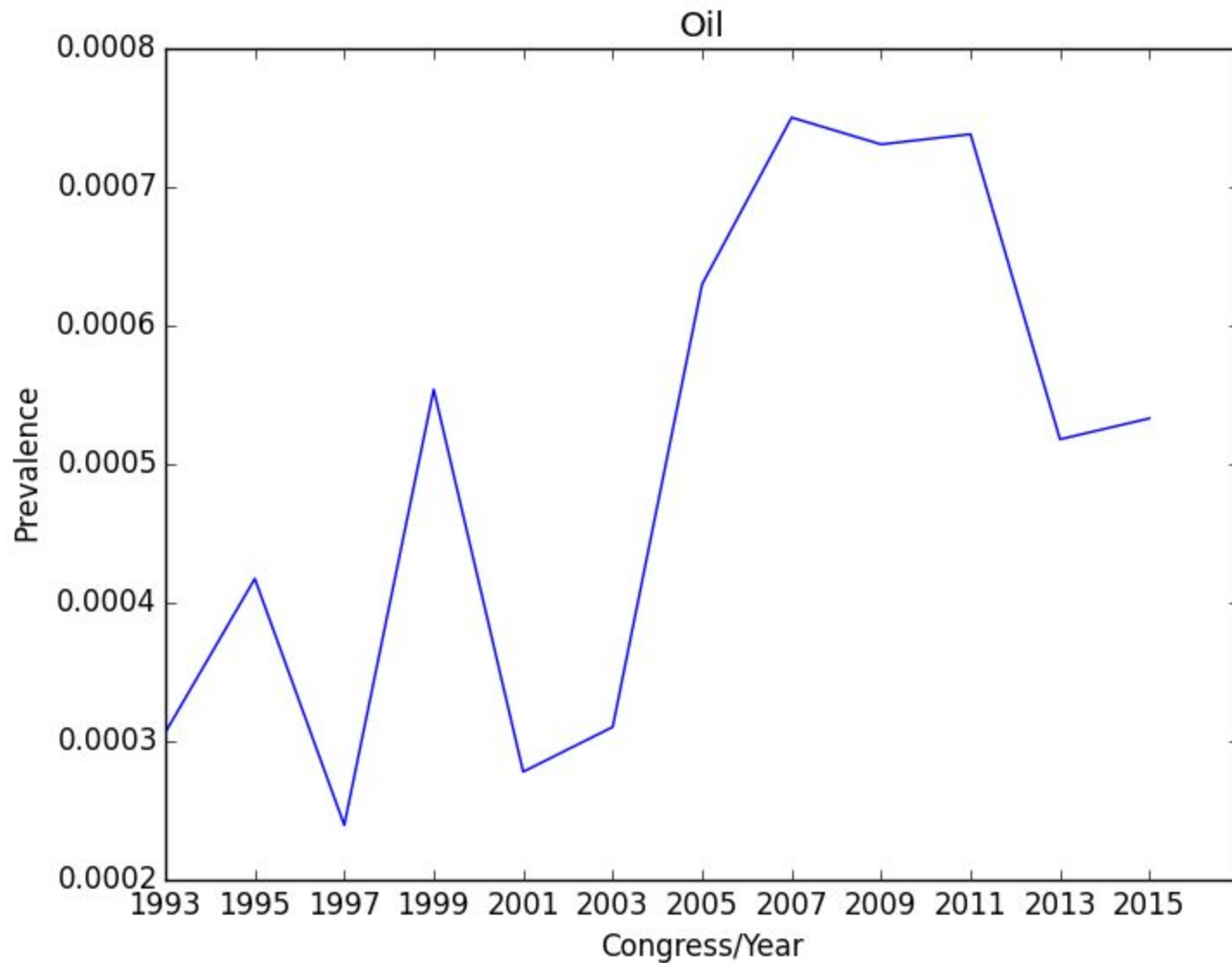
### *Affordable Care in Bills Over Time*



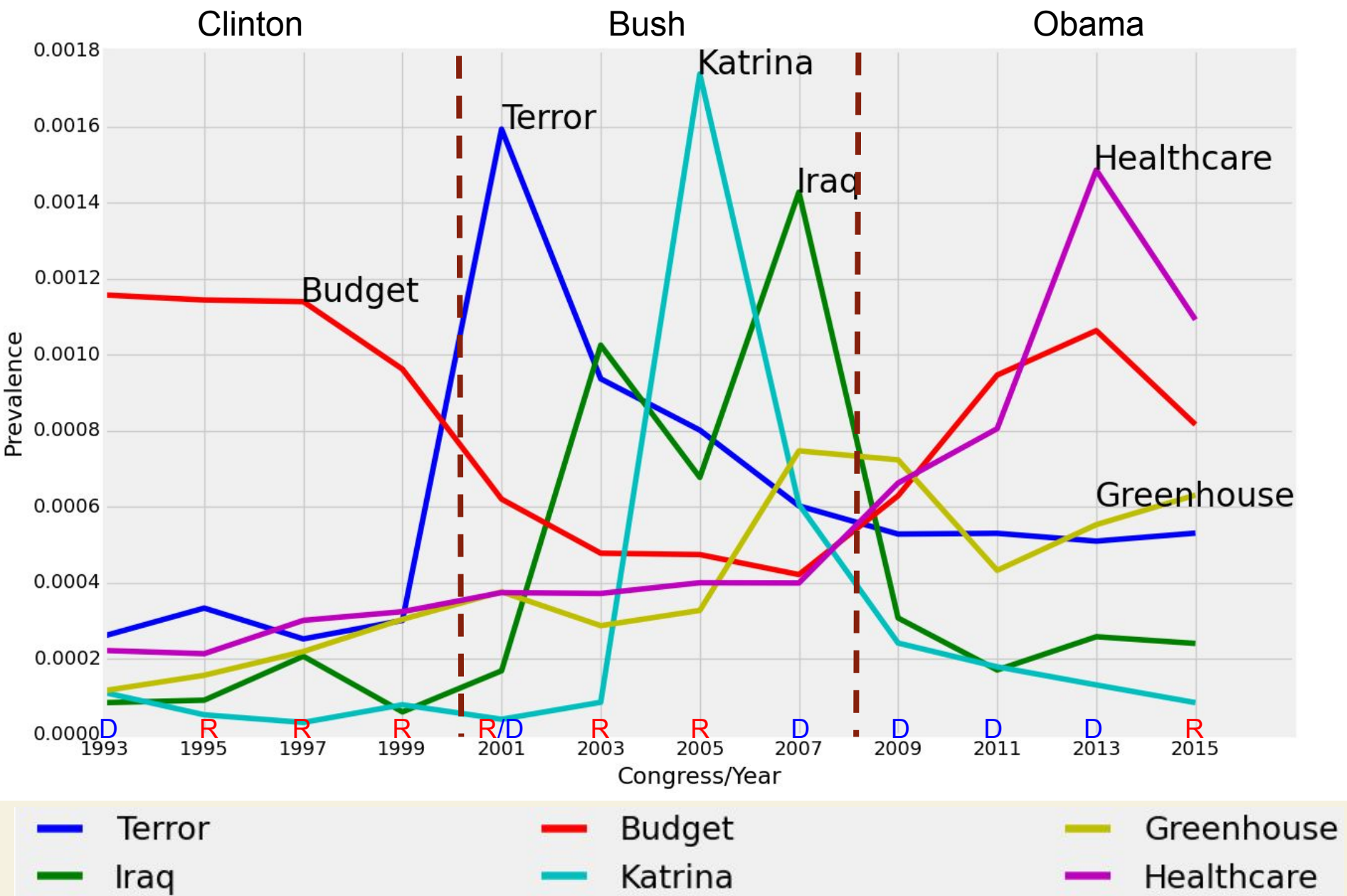
# Topic 235



## Oil



### *Oil in Bills Over Time*





Thanks!

**ANY QUESTIONS?**



You can find me at:

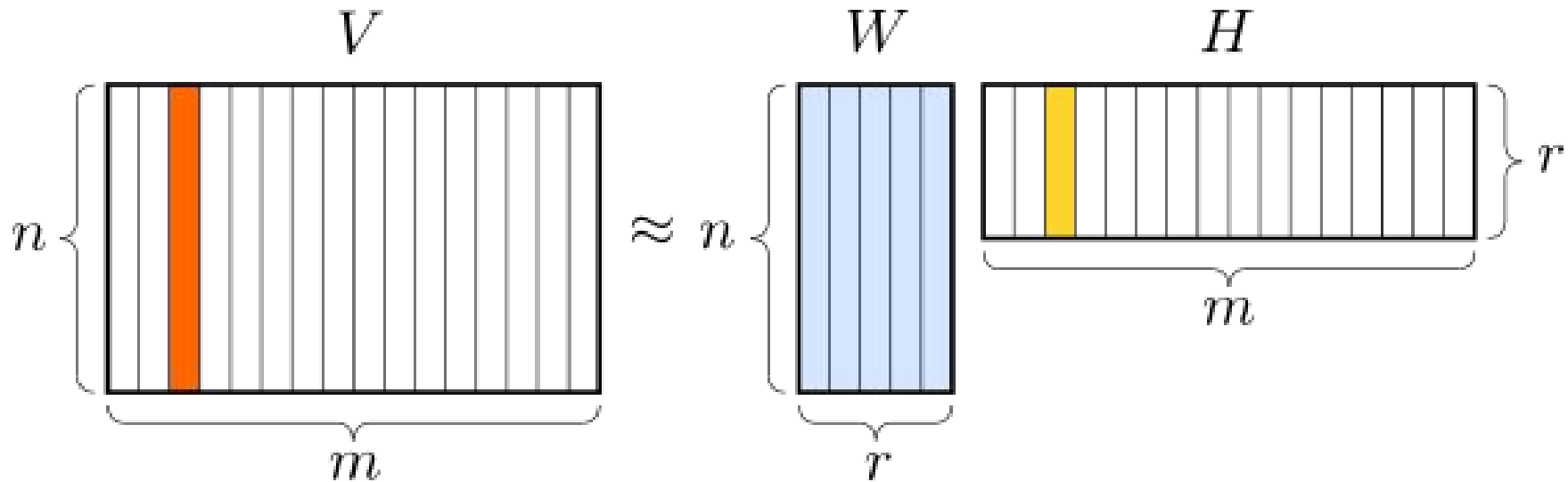
[linkedin.com/in/samuelcsherman](https://linkedin.com/in/samuelcsherman)

[github.com/scsherm](https://github.com/scsherm)

[scsherm@gmail.com](mailto:scsherm@gmail.com)

# Appendix

# Non-Negative Matrix Factorization

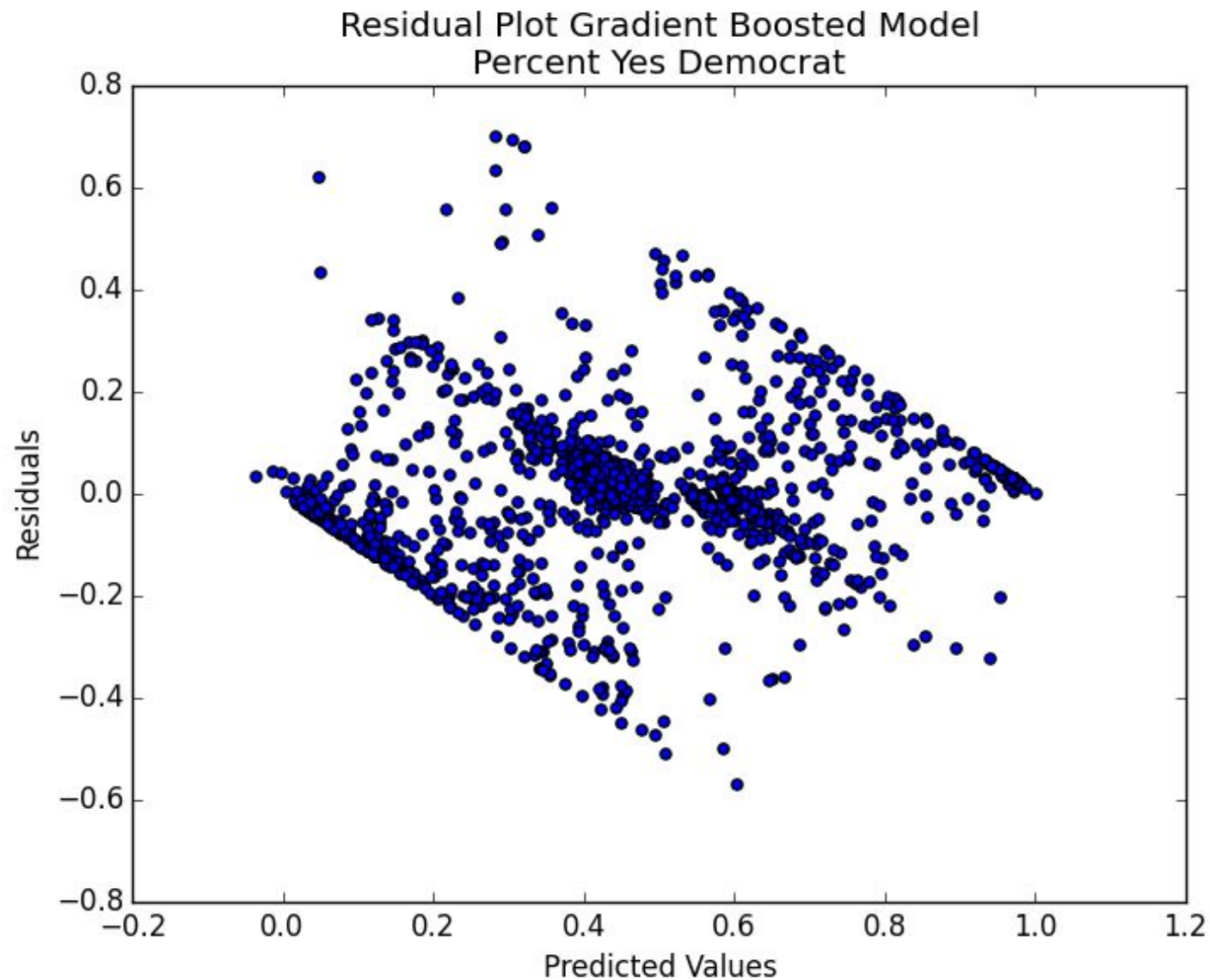


- Recommendation Systems
  - Null space
- Trend modeling/finance/stocks
  - Social Media
- Even Image Processing
  - Pattern/facial recognition

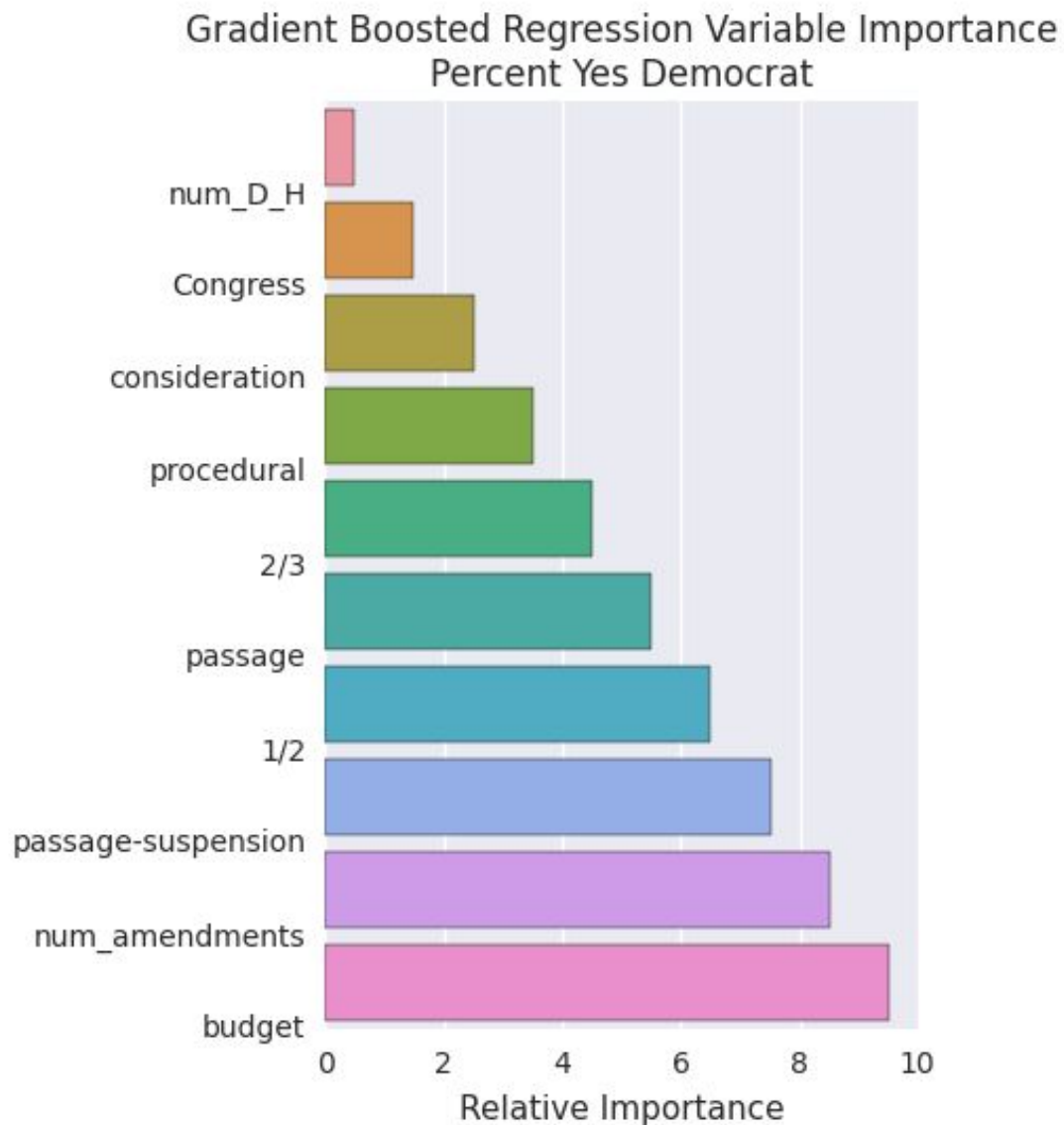
# Regression

## Percent Yes Democrat

	Mean Squared Error	Root Mean Squared Error	R <sup>2</sup> score
Random Forest	<b>0.0205</b>	<b>0.143</b>	<b>0.740</b>
Bagging	<b>0.0207</b>	<b>0.144</b>	<b>0.737</b>
Linear	<b>0.0417</b>	<b>0.204</b>	<b>0.417</b>
Gradient Boosted	<b>0.0203</b>	<b>0.143</b>	<b>0.743</b>



*Residuals*



## *Feature Importances*

# Classifier

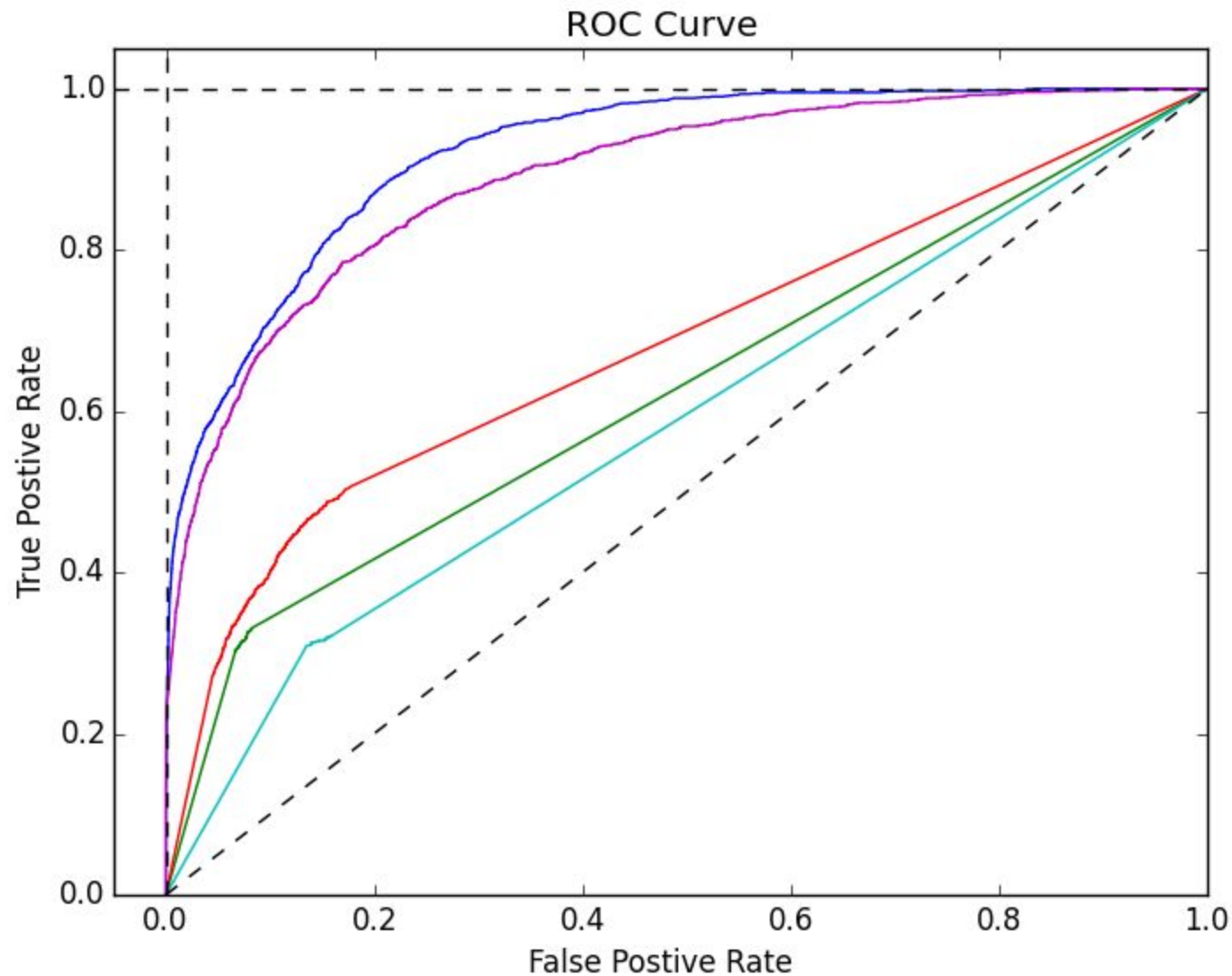


# 130,000 Bills

Since 1993

# 6% voted on

7,600

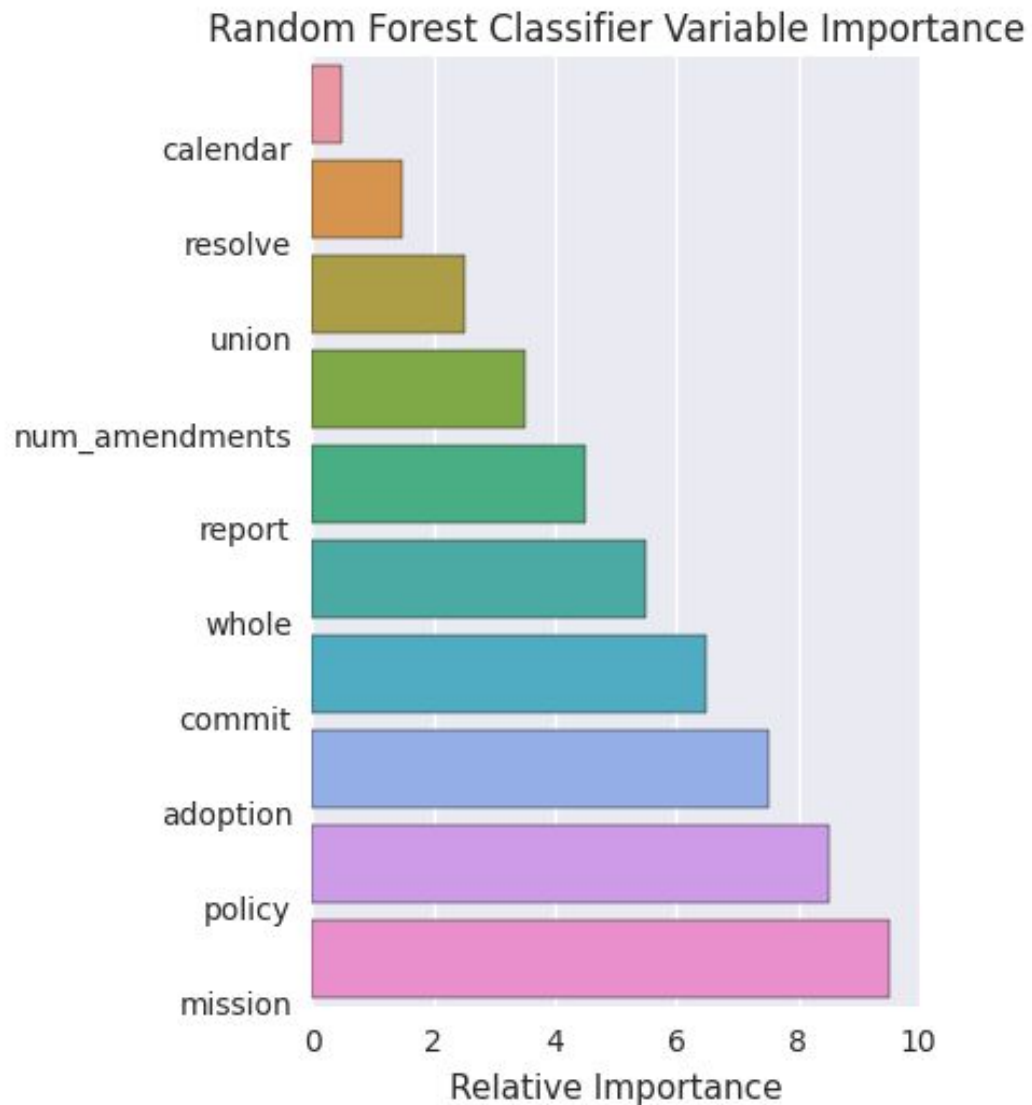


— RFC AUC = 0.921  
— GNB AUC = 0.626

— MNB AUC = 0.681  
— BNB AUC = 0.584

— logr AUC = 0.89

- ◉ **Oversampling**
  - **SMOTE**
- ◉ **Undersampling**
- ◉ **Recall/Precision**



## *Feature Importances*