



Module #9b:

Visualization – Text as Data



Visualization Techniques

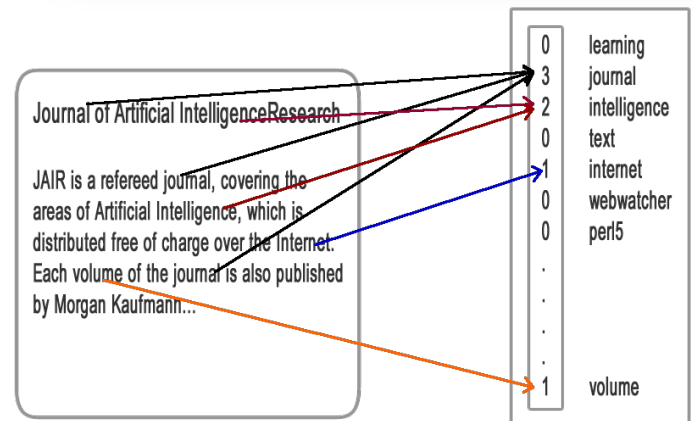
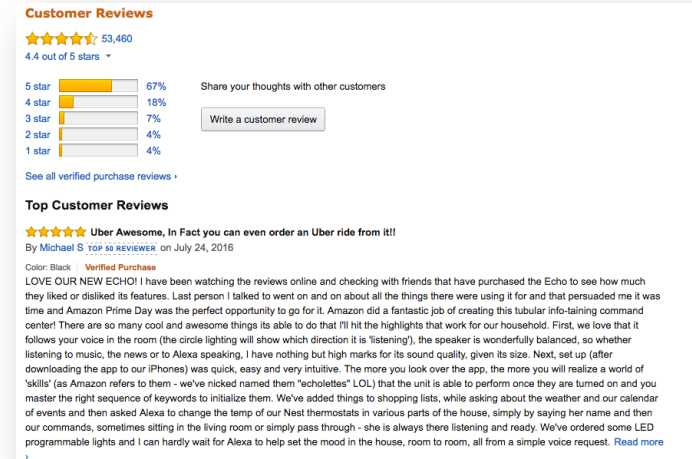
1. Text as Data
2. Visualizing Document Content
3. Evolving Documents
4. Visualizing Conversation
5. Document Collections



1. TEXT AS DATA

Bag of Words Model

- Ignore ordering relationships within the text
- A document \approx vector of term weights
- Each dimension corresponds to a term (10,000+)
- Each value represents the relevance (e.x. simple term counts)
- Aggregate into a document-term matrix
- Document vector space model





Word Count

WORDCOUNT

◀ PREVIOUS WORD

NEXT WORD ▶

the of and to ain that it is was i for on you he be with by have in the not but with is from which are done coming under all the words in the English language

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100

CURRENT WORD

FIND WORD:

BY RANK:

REQUESTED WORD: THE

RANK: 1

86800 WORDS IN ARCHIVE

[ABOUT WORDCOUNT](#)

<http://wordcount.org>

WordCounts (Harris '04)

Visualizations : Wordle of Sarah Palin RNC 9/3/2008 Speech

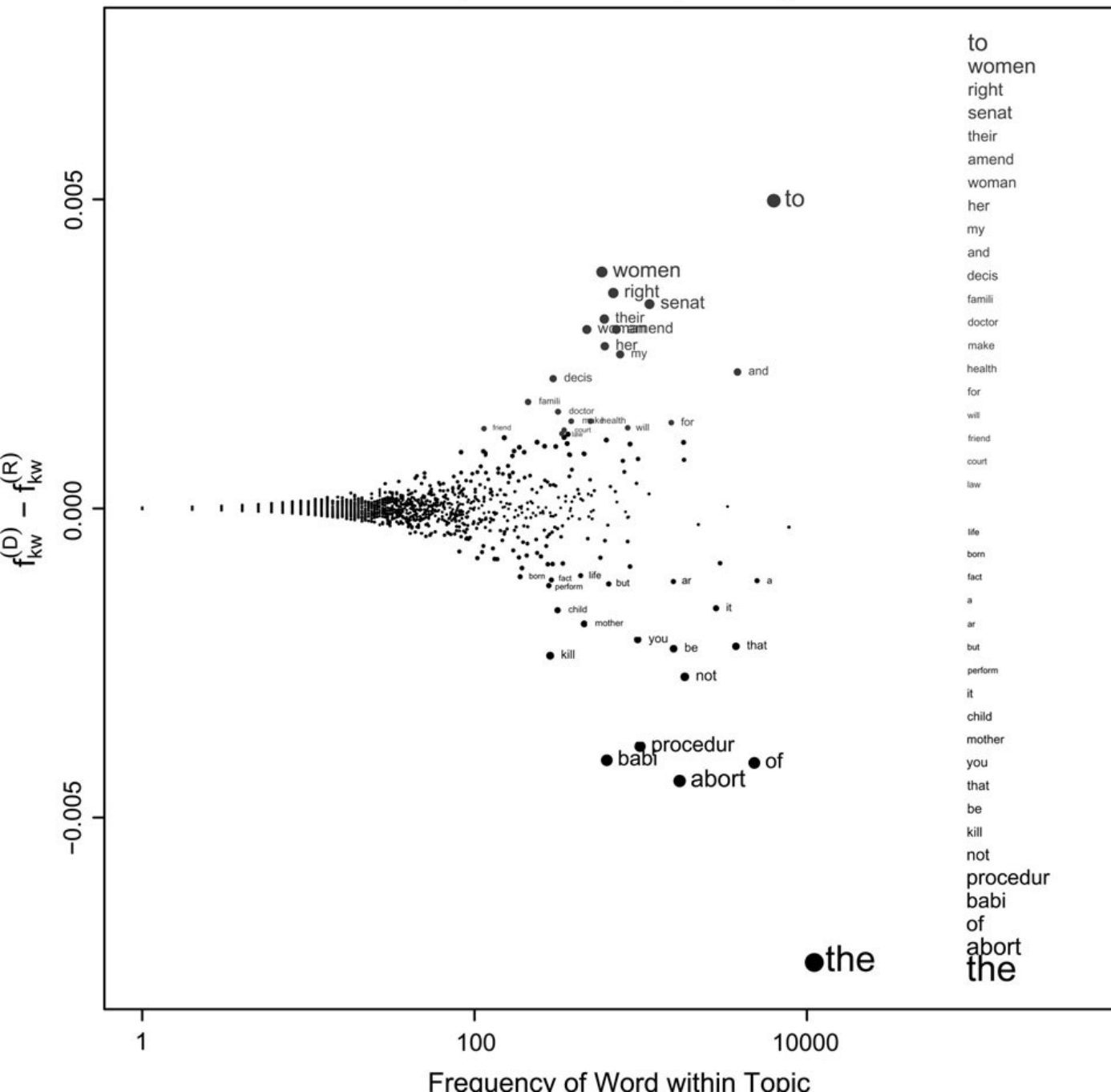
Creator: Anonymous

Tags:

Edit Language Font Layout Color



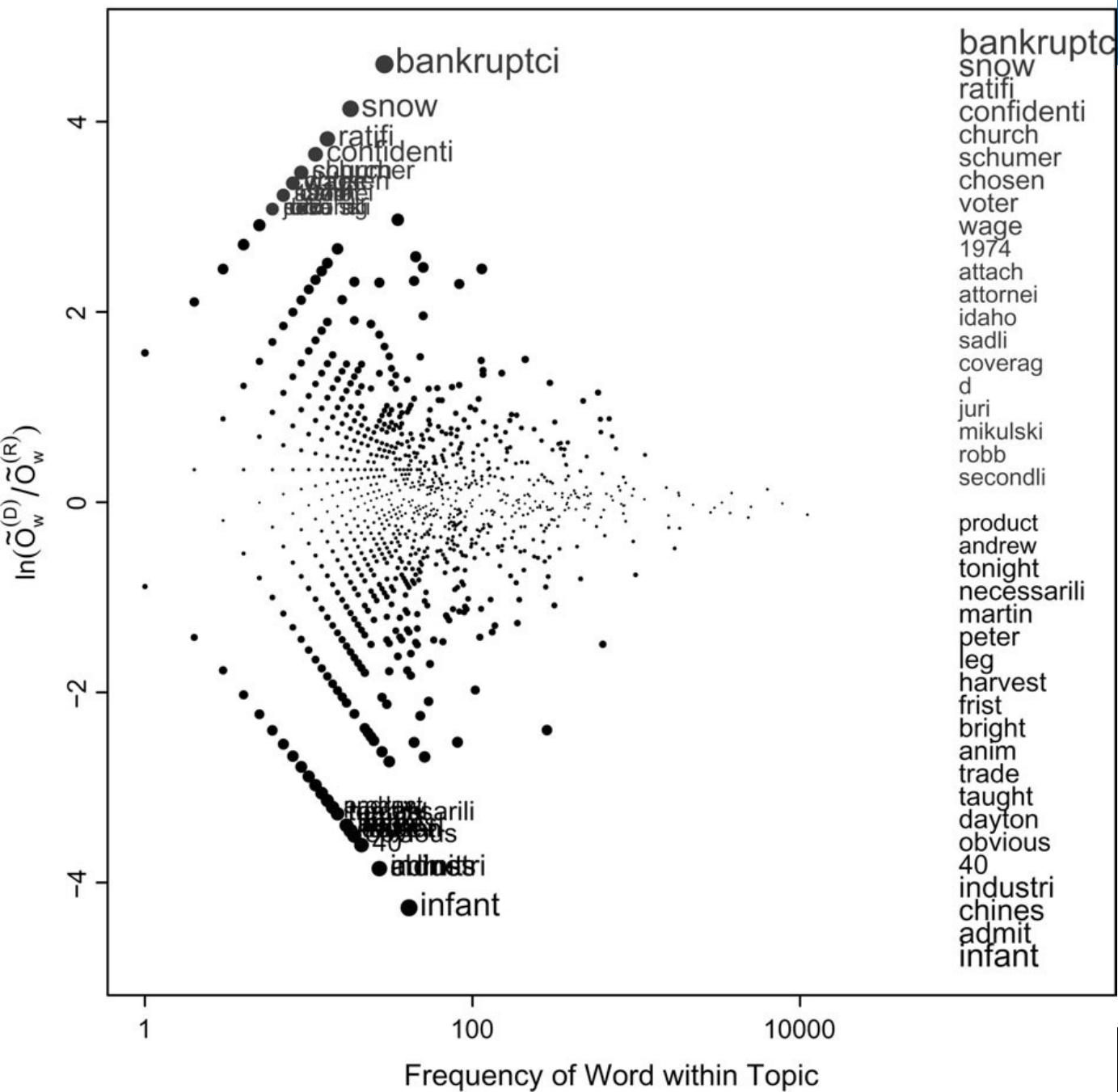
Partisan Words, 106th Congress, Abortion (Difference of Proportions)



Senate speech from
1997 to 2004

Burt L. Monroe et al. ' Words:
Lexical Feature Selection and
Evaluation for Identifying the
Content of Political Conflict.
Polit Anal 2008; 16 (4): 372-403.

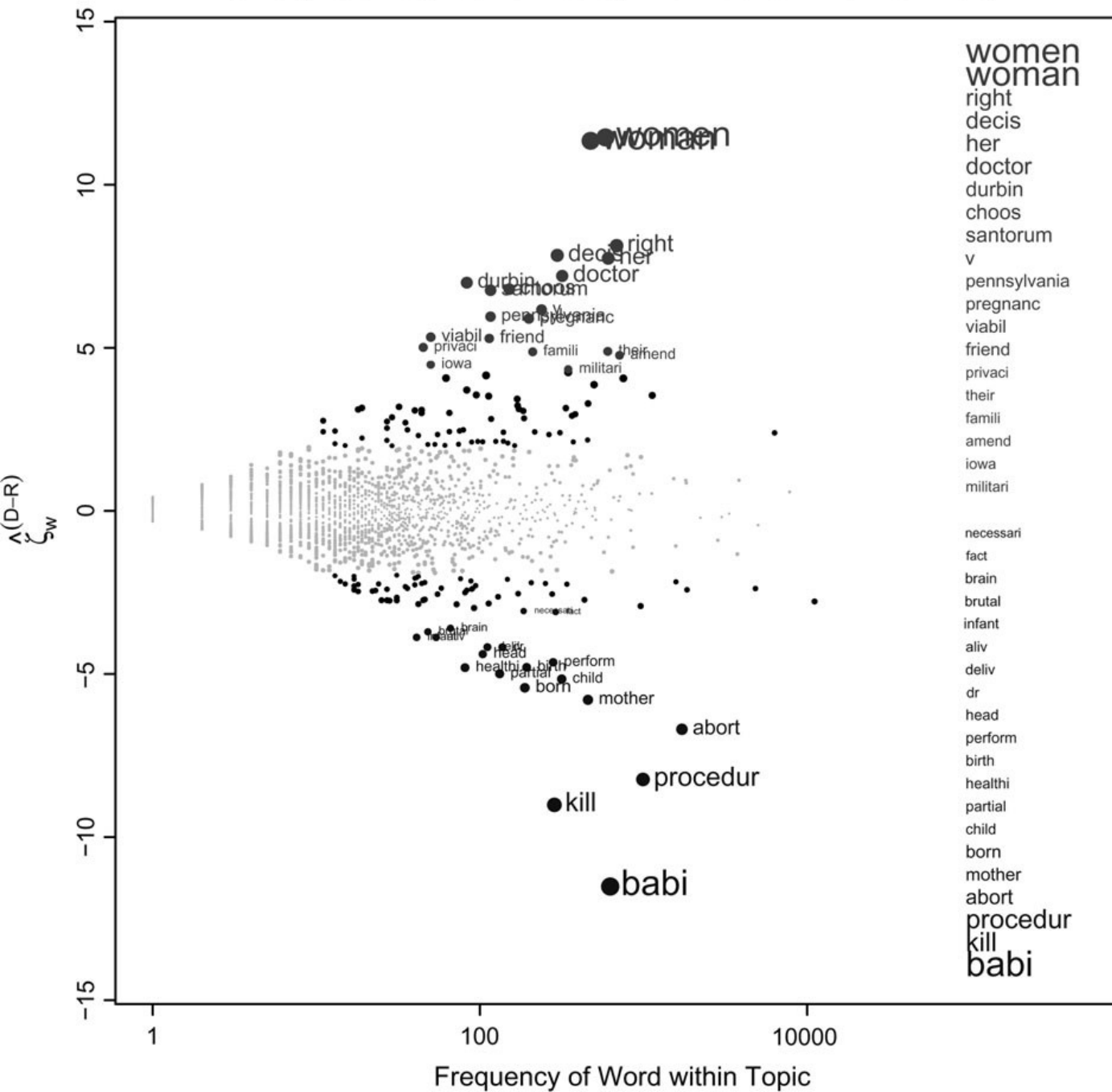
Partisan Words, 106th Congress, Abortion (Log-Odds-Ratio, Smoothed Log-Odds-Ratio)



Senate speech from
1997 to 2004

Burt L. Monroe et al. ' Words:
Lexical Feature Selection and
Evaluation for Identifying the
Content of Political Conflict.
Polit Anal 2008; 16 (4): 372-403.

Partisan Words, 106th Congress, Abortion (Weighted Log-Odds-Ratio, Informative Dirichlet Prior)



Senate speech from
1997 to 2004

Burt L. Monroe et al. ' Words:
Lexical Feature Selection and
Evaluation for Identifying the
Content of Political Conflict.
Polit Anal 2008; 16 (4): 372-403.



Keyword Weighting

Term Frequency

- $\text{tftd} = \text{count}(t) \text{ in } d$
- Can take log frequency: $\log(1 + \text{tftd})$
- Can normalize to show proportion: $\text{tftd} / \sum_t \text{tftd}$

Term frequency–inverse document frequency (tf-idf)

- Reduce the weight of terms that occur very frequently in the document(s) and increase the weight of terms that occur rarely.
- Example: the vs. healthcare

- TF.IDF: Term Freq by Inverse Document Freq

$$\text{tf.idftd} = \log(1 + \text{tftd}) \times \log(N/\text{dft})$$

$\text{dft} = \# \text{ docs containing } t;$

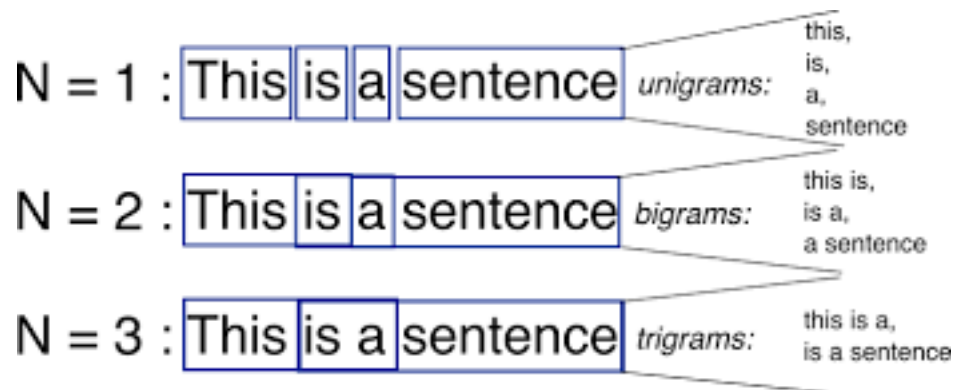
$N = \# \text{ of docs}$

Limitations of Freq. Statistics

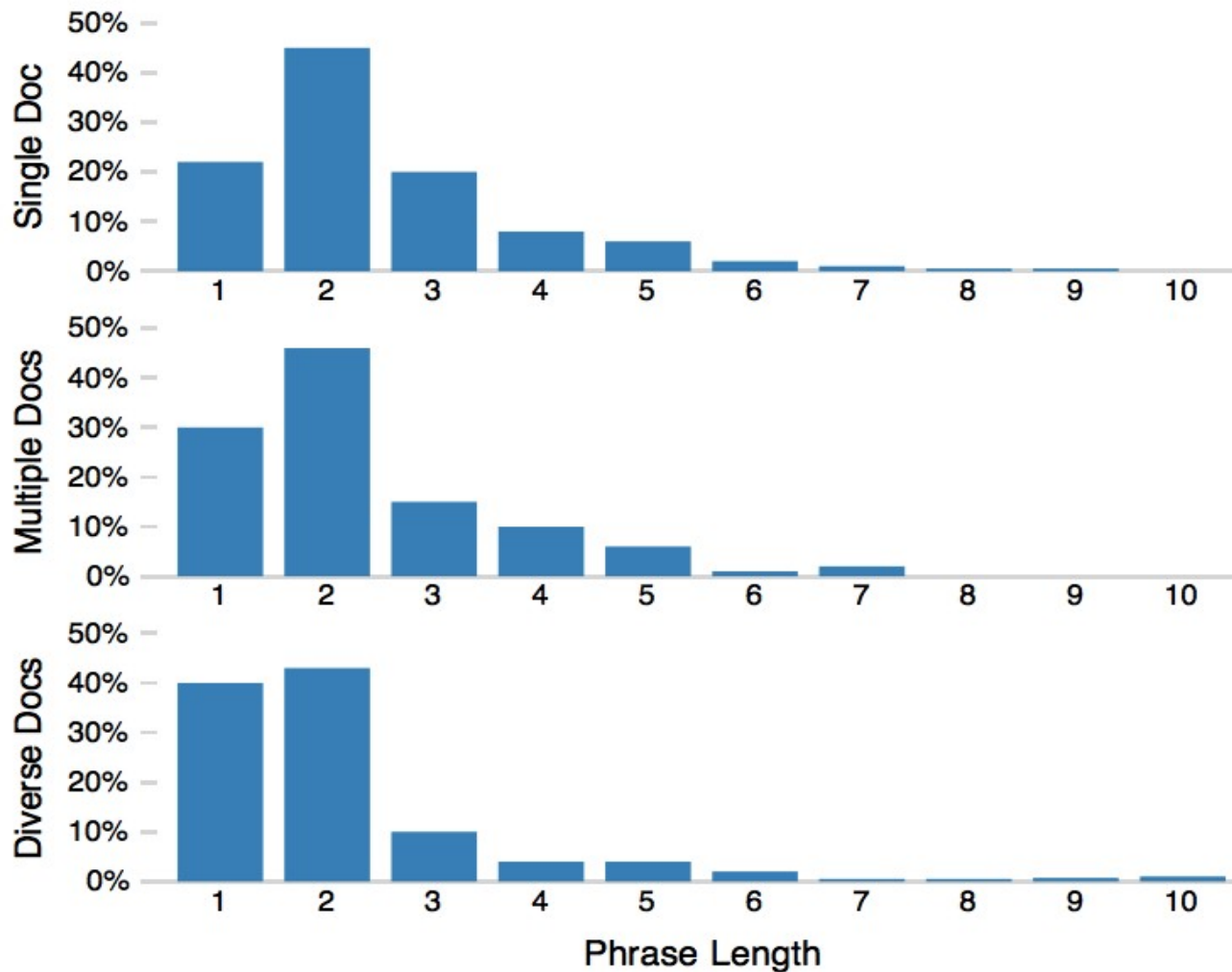
- Typically focus on unigrams (single terms)
- Often favors frequent (TF) or rare (IDF) terms
 - Not clear that these provide best description



- A “bag of words” ignores additional information
 - Grammar / part-of-speech $N =$
 - Position within document
 - Recognizable entities $N =$



N-grams



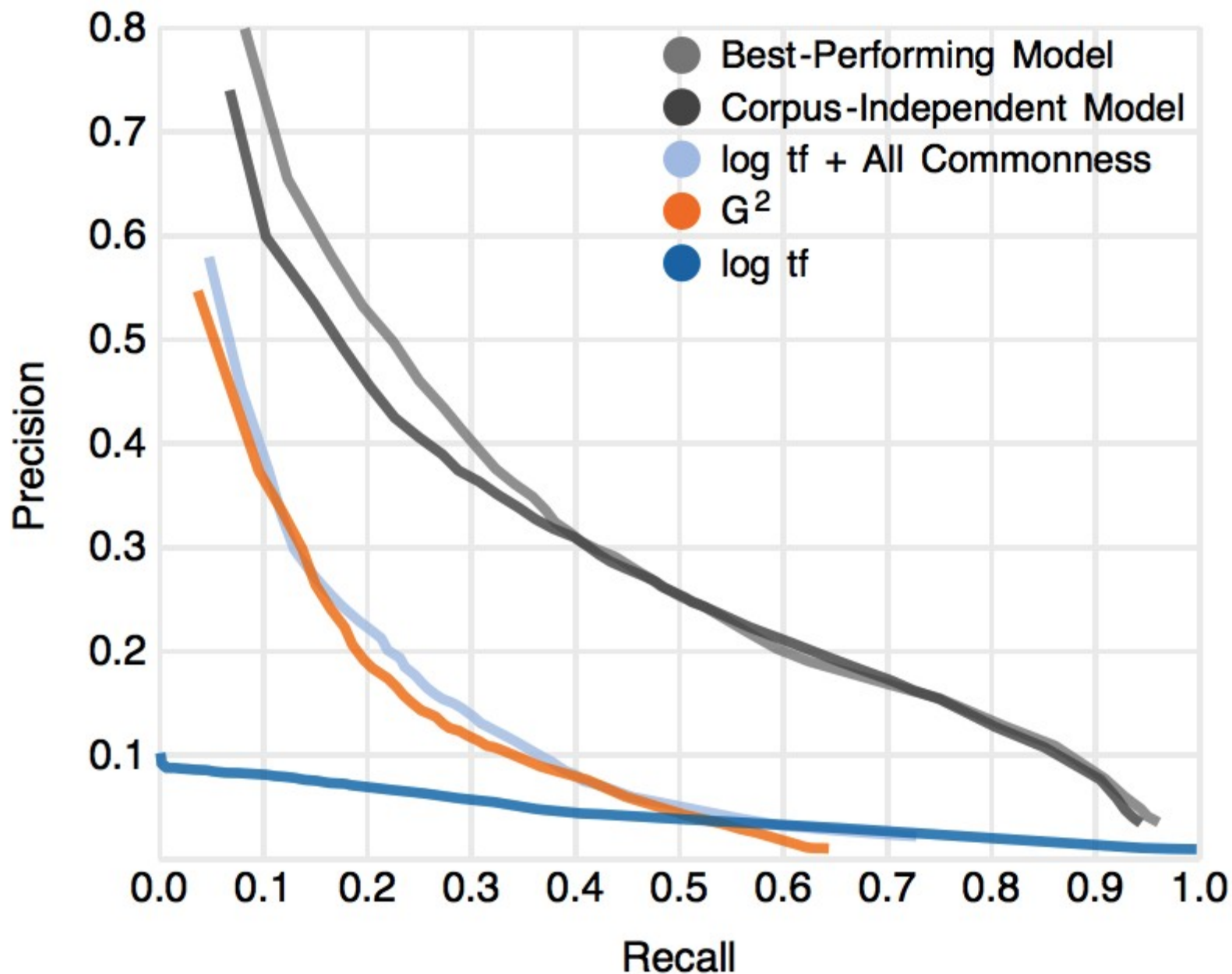


Term Commonness

$$\log(\text{tf}_w) / \log(\text{tf}_{\text{the}})$$

- The normalized term frequency relative to the most frequent n-gram, e.g., the word “the”.
- Measured across a corpus or across the entire English language (using Google n-grams)

Scoring Terms with Freq, Grammar & Position





A fighter jet rain check

Story and video by [Chamila Jayaweera](#)

Have you ever thought about what it takes to make sure that sea-based fighter jets stay dry?

When it comes to the F/A-18 Super Hornet, Boeing engineers in St. Louis use a special process called the Water Check Test to rule out areas where moisture could seep into the aircraft and its electronics suite.

Program experts douse the jet with simulated rain at a 15-inch-per-hour rate for about 20 minutes inside an enormous hangar in St. Louis.

"Our ultimate customers are U.S. Navy fighter pilots, and we want to ensure their safety in flight and on the ground, and water-tight integrity of the aircraft also helps increase their effectiveness," said Boeing's Rich Baxter, F/A-18 Super Hornet final assembly manager.

To find out more about how the process works and watch the action unfold, click above to see the video story.



CHAMILA JAYAWEERA/BOEING

The Water Check team rolls in a large metal frame, which they affectionately call their "spray tree," over a Super Hornet inside a St. Louis hangar.



G2

Regression Model

fighter

F/A

Hornet

Super

Boeing

-18

rain

St.

jet

Louis

15-inch-per-hour

douse

hangar

water-tight

Check

Baxter

sea-based

aircraft

Rich

seep

click

Navy

sure

Water

moisture

watch

enormous

stay

want

Super Hornet

F/A -18

fighter jet

Boeing engineers

special process

rain check

electronics suite

Program experts

simulated rain

ultimate customers

enormous hangar

water-tight integrity

Rich Baxter

15-inch-per-hour rate

video story

aircraft

U.S. Navy fighter pilots

Super Hornet final assembly manager

U.S.
Navy fighter
fighter pilot
sea-based fighter



Tips for Effective Test-As-Data Visualizations

- Understand the limitations of your language model.
- Bag of words:
 - Easy to compute
 - Single words
 - Loss of word ordering
- Select appropriate model and visualization
- Generate longer, more meaningful phrases Adjective-noun word pairs for reviews
- Show key phrases within source text