



# BEST SELLING BOOKS

What topics are people reading?



# OBJECTIVE: CATEGORIZE POPULAR BOOKS SINCE 2008

---

## BOOKS

# The New York Times Best Sellers

Authoritatively ranked lists of books sold in the United States, sorted by format and genre.

✓ FICTION

Combined Print & E-Book Fiction

**Hardcover Fiction**

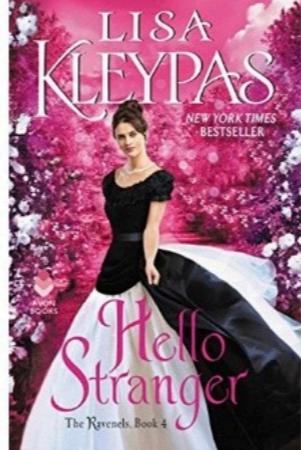
Paperback Trade Fiction

Audio Fiction

EBOOKS | MONTHLY LISTS

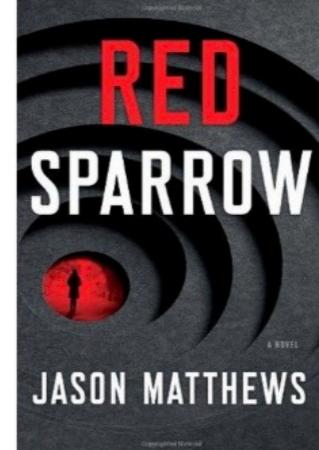
Fiction >

2



LISA KLEYPAS  
NEW YORK TIMES BESTSELLER  
*Hello Stranger*  
The Ravenels, Book 4

3



RED SPARROW  
JASON MATTHEWS

4 WEEKS ON THE LIST

**THE GREAT ALONE**

by Kristin Hannah

A former prisoner of war returns from Vietnam and moves his family to Alaska, where they face tough conditions.

NEW THIS WEEK

**HELLO STRANGER**

by Lisa Kleypas

Book 4 of the Ravenels series.

2 WEEKS ON THE LIST

**RED SPARROW**

by Jason Matthews

A Russian intelligence officer trained in the art of seduction becomes entangled with a young C.I.A. officer.

# BUSINESS CASE & HYPOTHESIS

---

## Potential Uses:

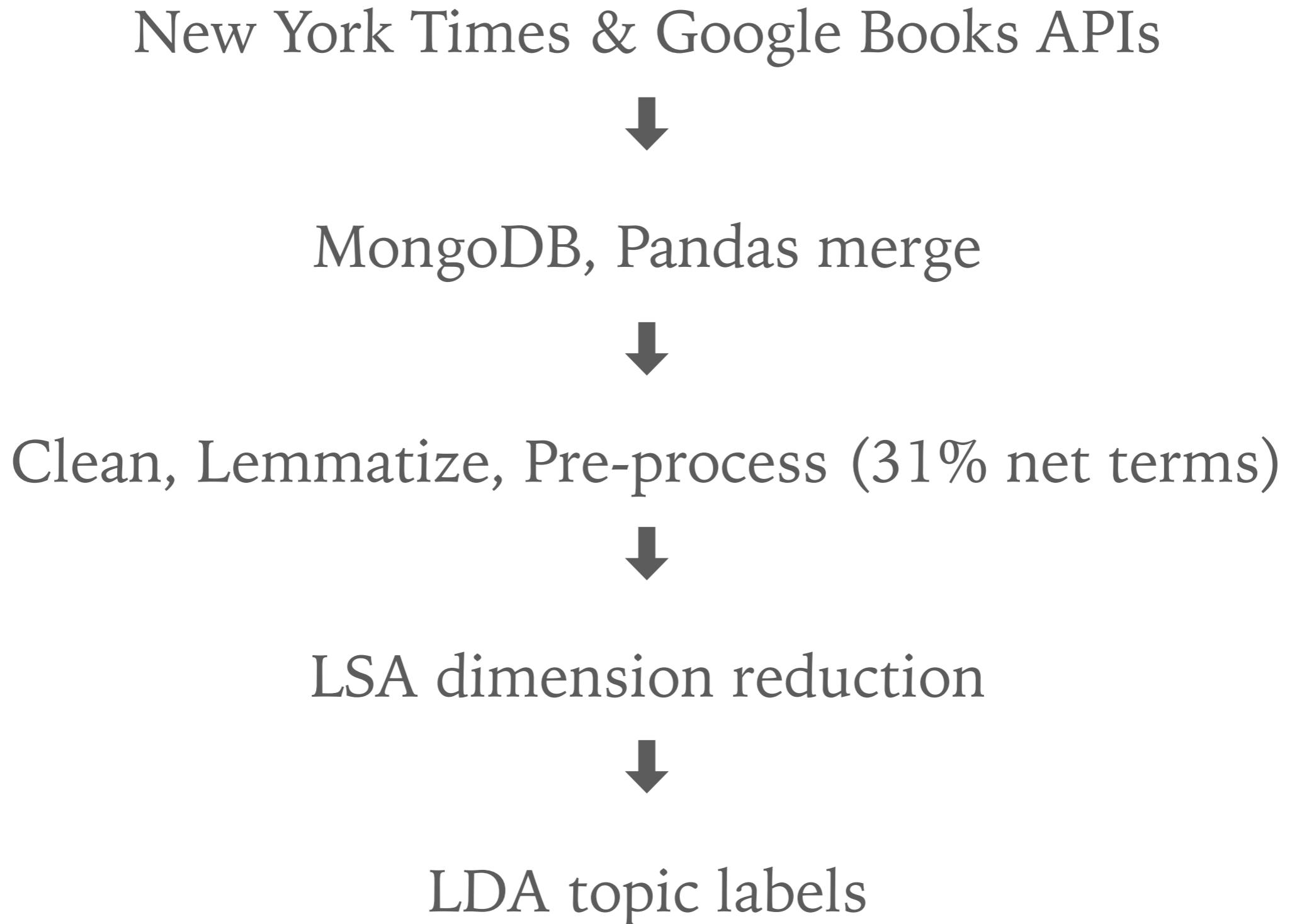
- Identify demand for authors/publishers
- Content-based recommendation for readers

## Hypothesis:

- We can find what people are most interested in by reducing the number of categories
  - 20+ for NYT articles
  - 30+ for Amazon books

# METHODOLOGY

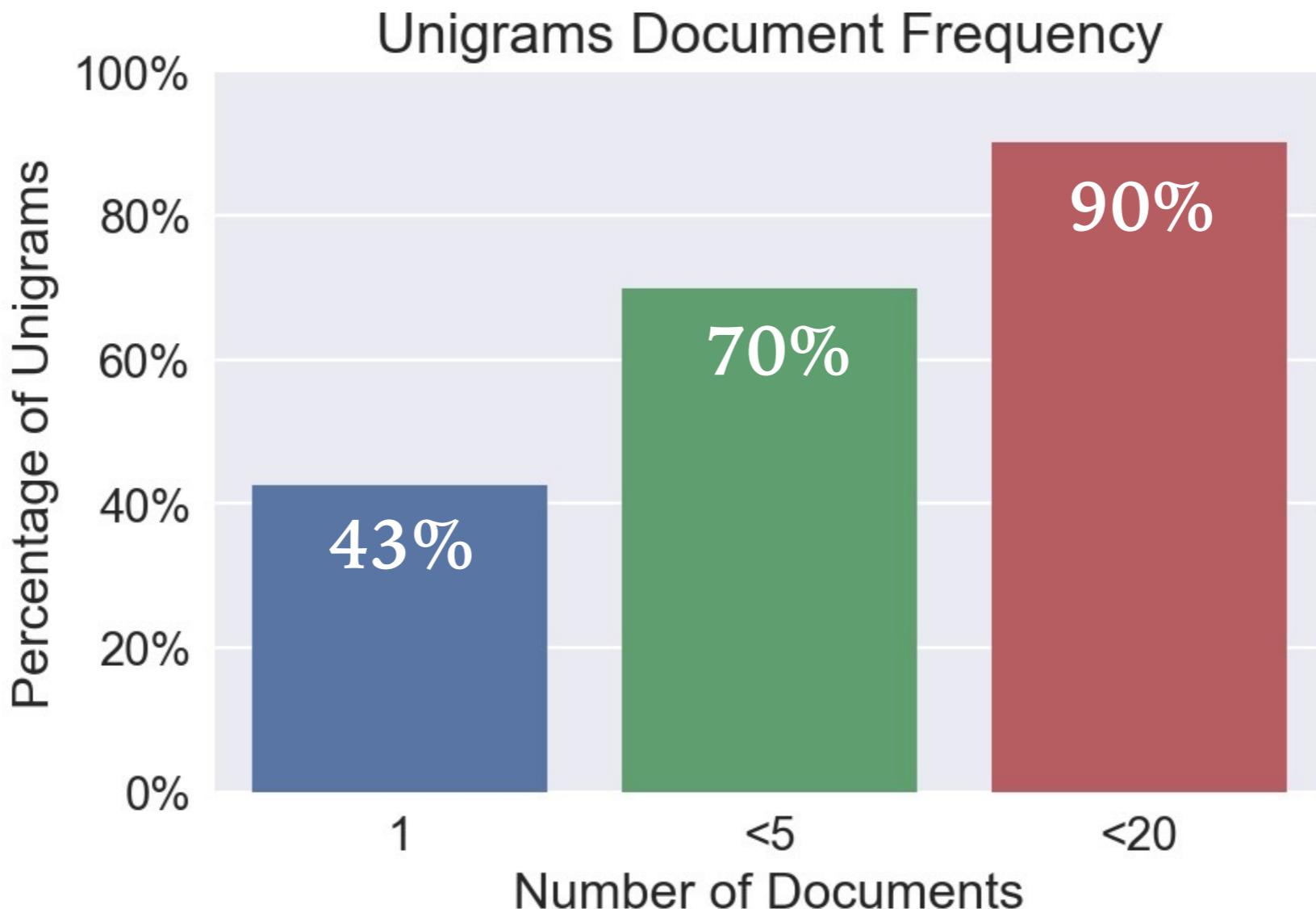
---



# 11,000 BOOKS: MOSTLY LONG TAIL N-GRAMS

---

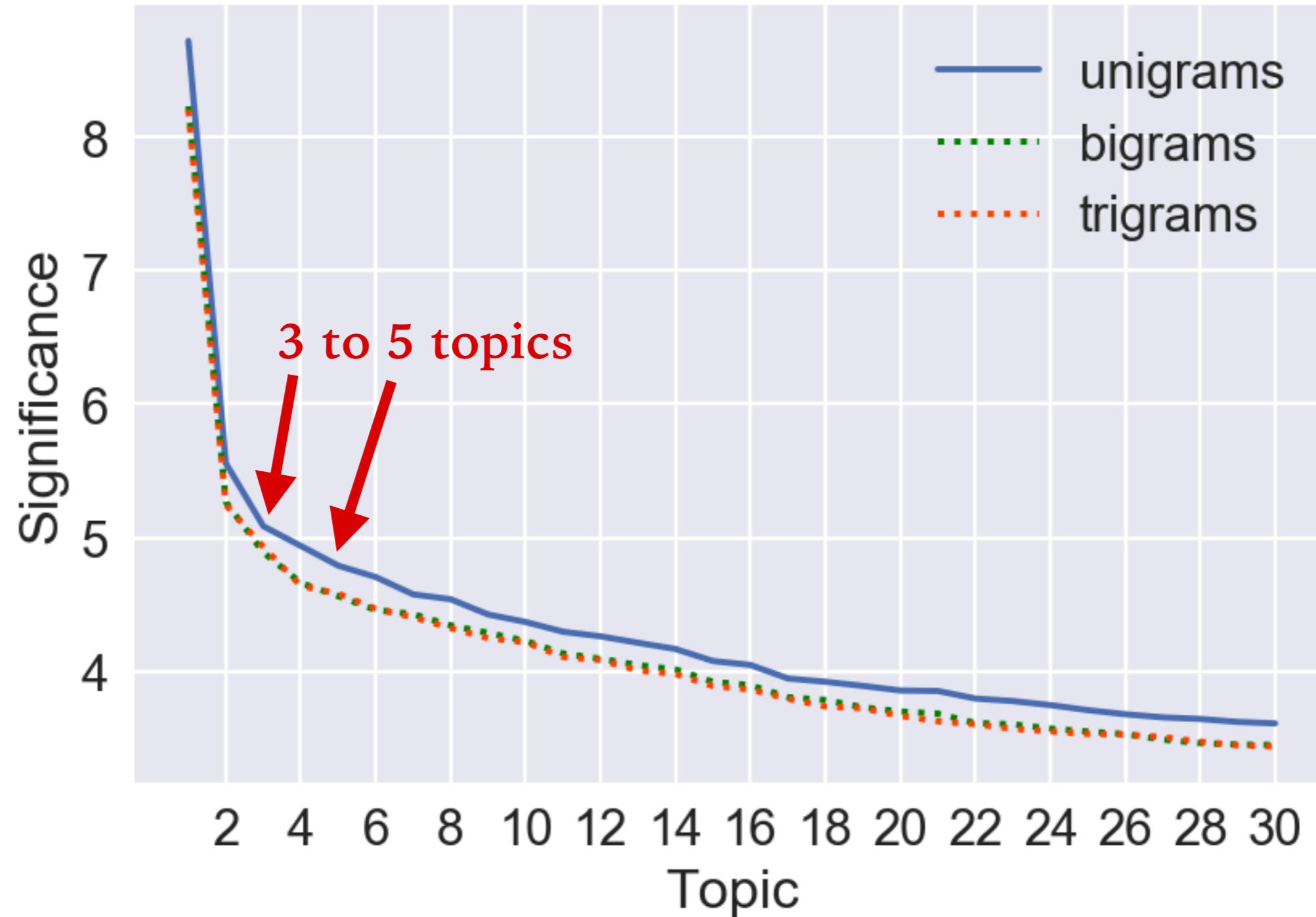
- Only 3 words occur in more than 10% of documents
- 20 documents represent 0.1% of corpus



# OPTIMAL DIMENSIONALITY

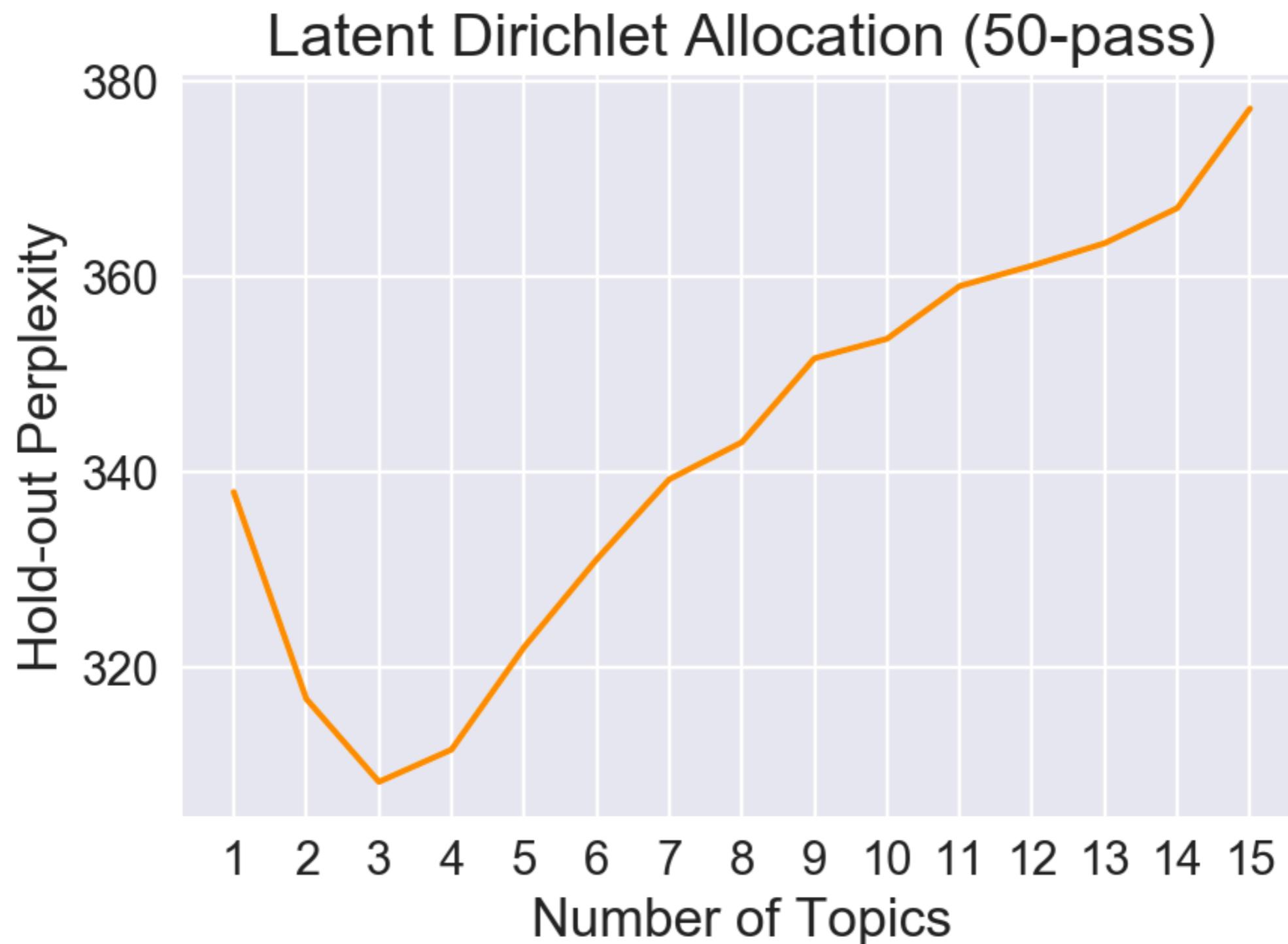
---

## Latent Semantic Analysis



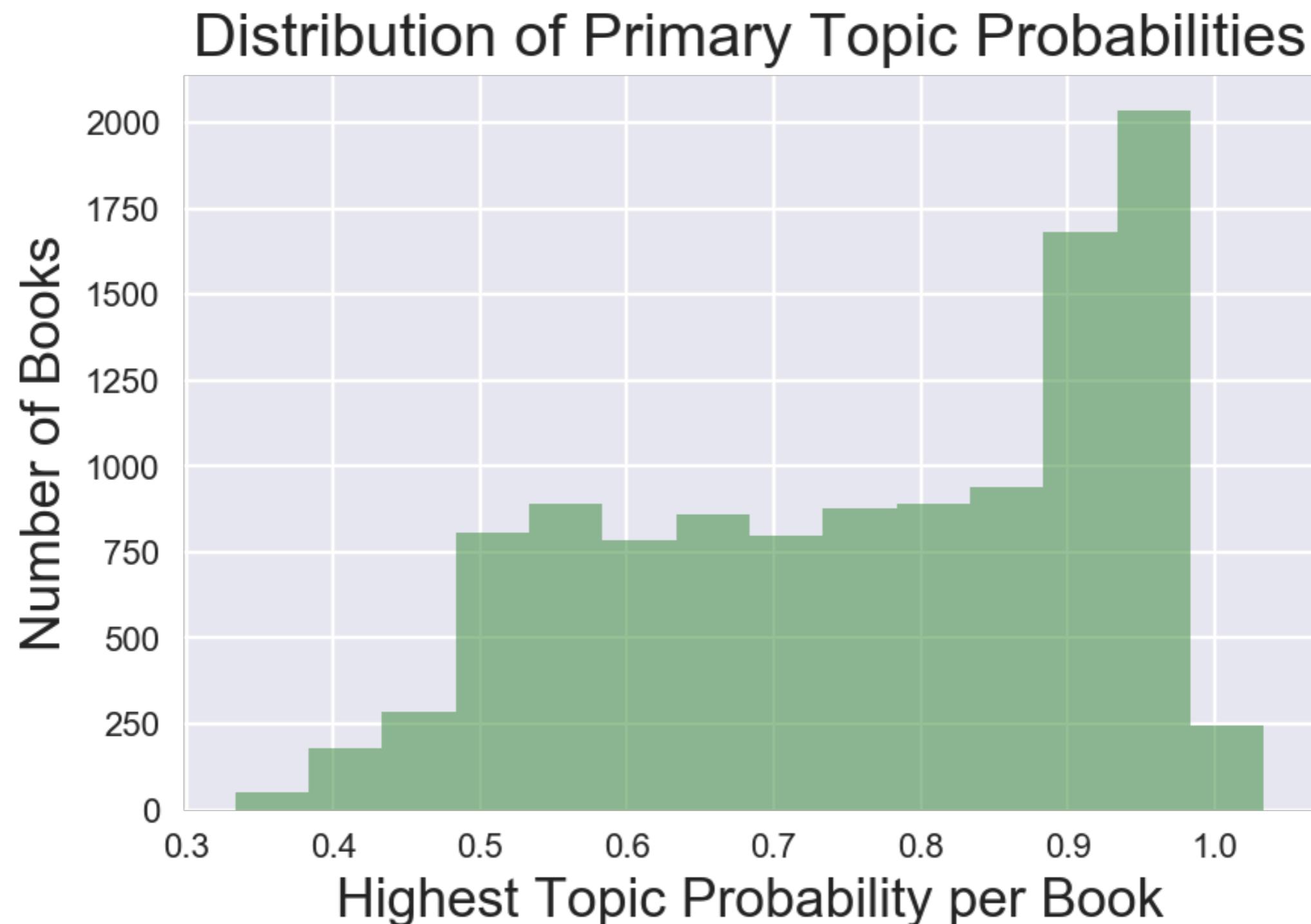
# OPTIMAL DIMENSIONALITY

---



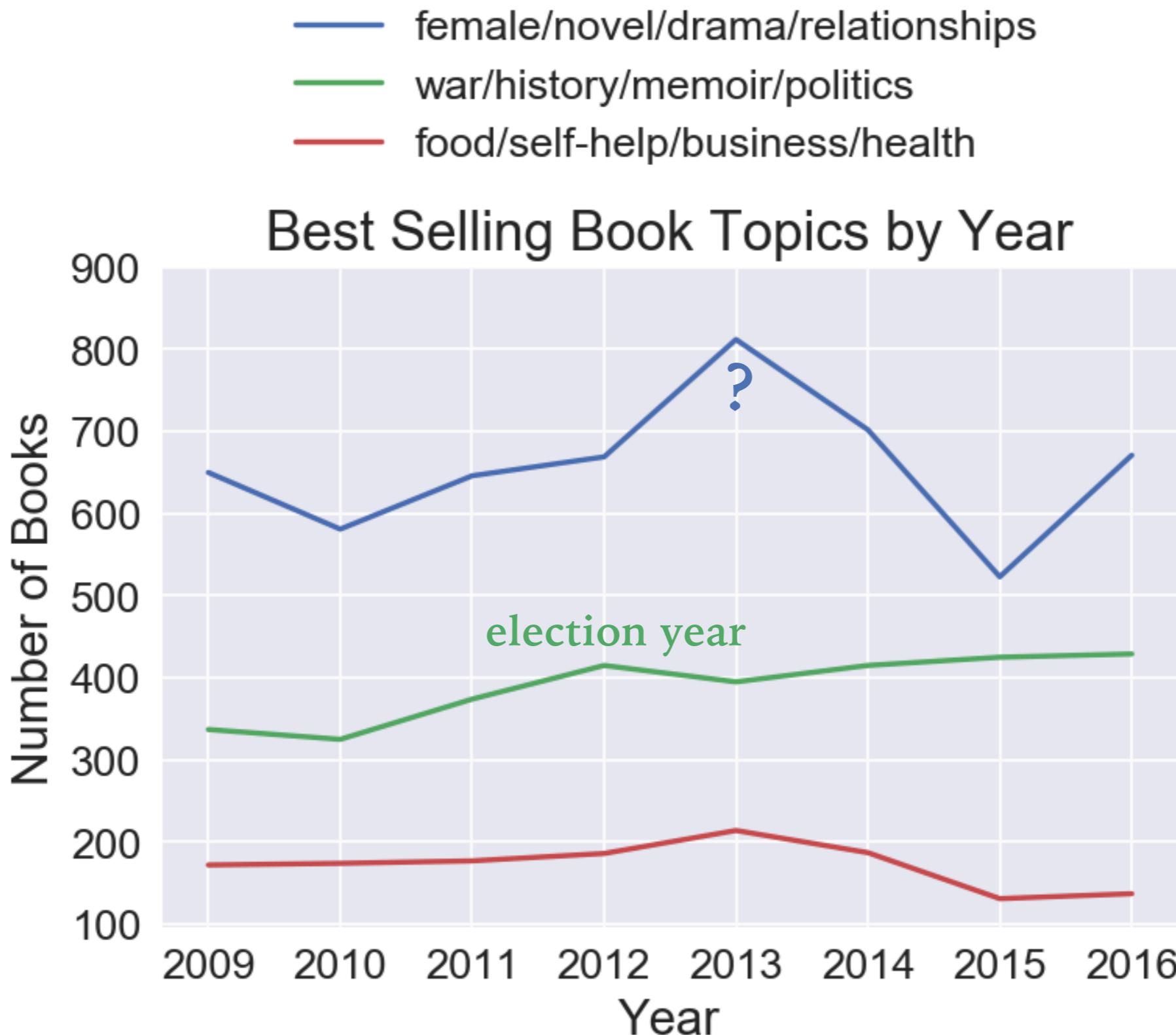
# LDA TOPIC PROBABILITIES

---



# LDA TOPIC TRENDS

---



# TOPIC QUALITY & ISSUES

.....

- Topics make sense, but limited use
- Term overlap, need more granularity
- Long tail n-grams, signal vs. noise
- Small dataset with mysterious ranking criteria
- Algorithm selection



A photograph of a stack of books. In the foreground, an open book is visible, showing its pages fanned out. The background consists of several more books stacked vertically, with their spines and parts of their pages visible.

*Thank you*