

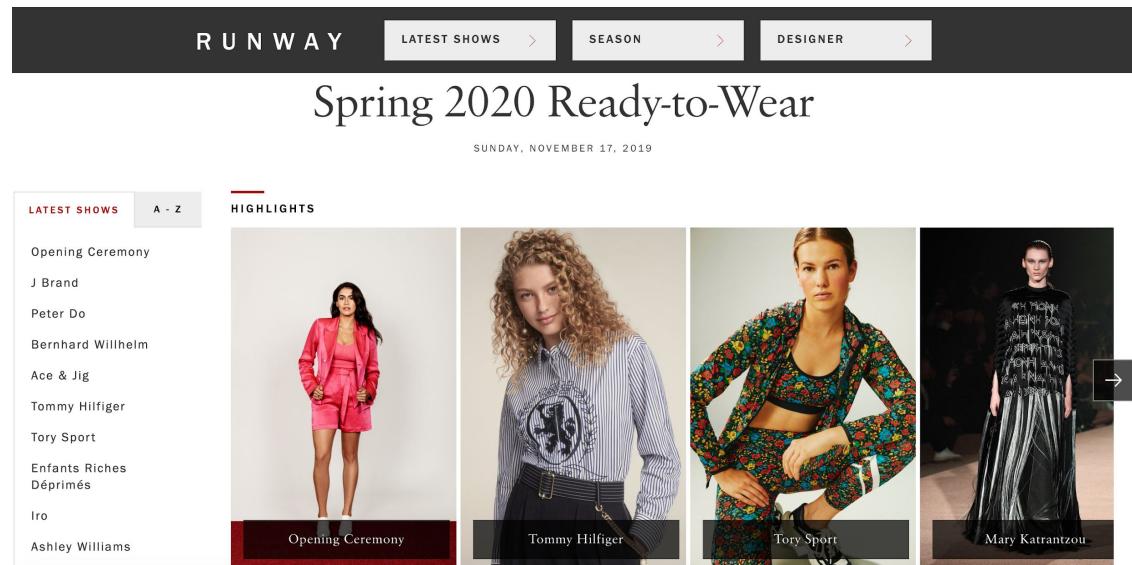
# RE-CLASSIFYING VOGUE COLLECTION COVERAGE



# Goal

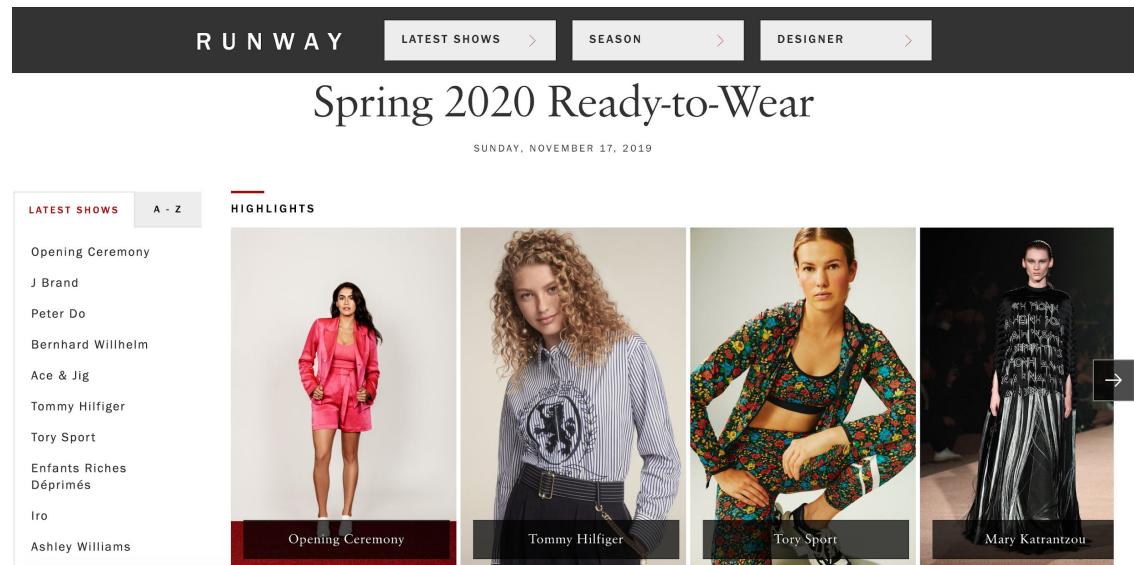
To reclassify Vogue runway show coverage from the Vogue website.

- ❑ What other categories can be created?
- ❑ How is the text for each collection related?
- ❑ Articles currently only categorized by season, year, and designer



# The Data

- ❑ Gathered 3525 articles covering major collections from 2018 to 2020
- ❑ Used Selenium to scrape from the Vogue website



# Bag of Words

- ❑ Created a bag of words.
  - ❑ These are the most frequently used words in the corpus.



# Once Common Stop Words Removed



# Further Cleaning

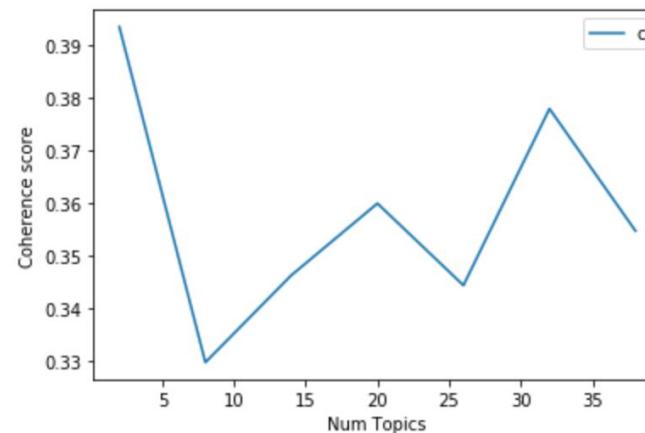
- Created a list of 100 most common words in corpus
- Used domain knowledge to pick terms (very general) to remove from corpus that would not help in modeling more specific topics (collection type, designer, year, and common terms like ‘runway’, ‘designer’, ‘collection’)
- Removed numbers and tokens with a length less than 3

# Next Steps

- ❑ Lemmatized the updated corpus
- ❑ Created bi-gram and tri-gram models with a minimum count of 10 and converted them into tokens
- ❑ Filtered extremes (since trying to find more specific niche words tokens must appear 10 times and appear in no more than 10% of the documents.)

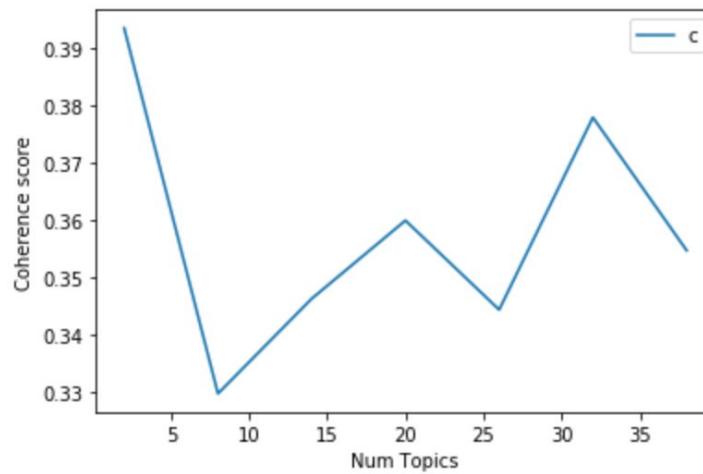
# First Model

- Created an LDA model to extract topic from the textual data.
- First chosen topic count = 5
- Coherence score was .374
- To improve model created a line graph showing how coherence scores increase by topic



# Choosing number of Topics

- ❑ First created 20 topics according to first peak
- ❑ Coherence score was .0.402
  - Improved by 5.5%
- ❑ Thought this might be too many topics



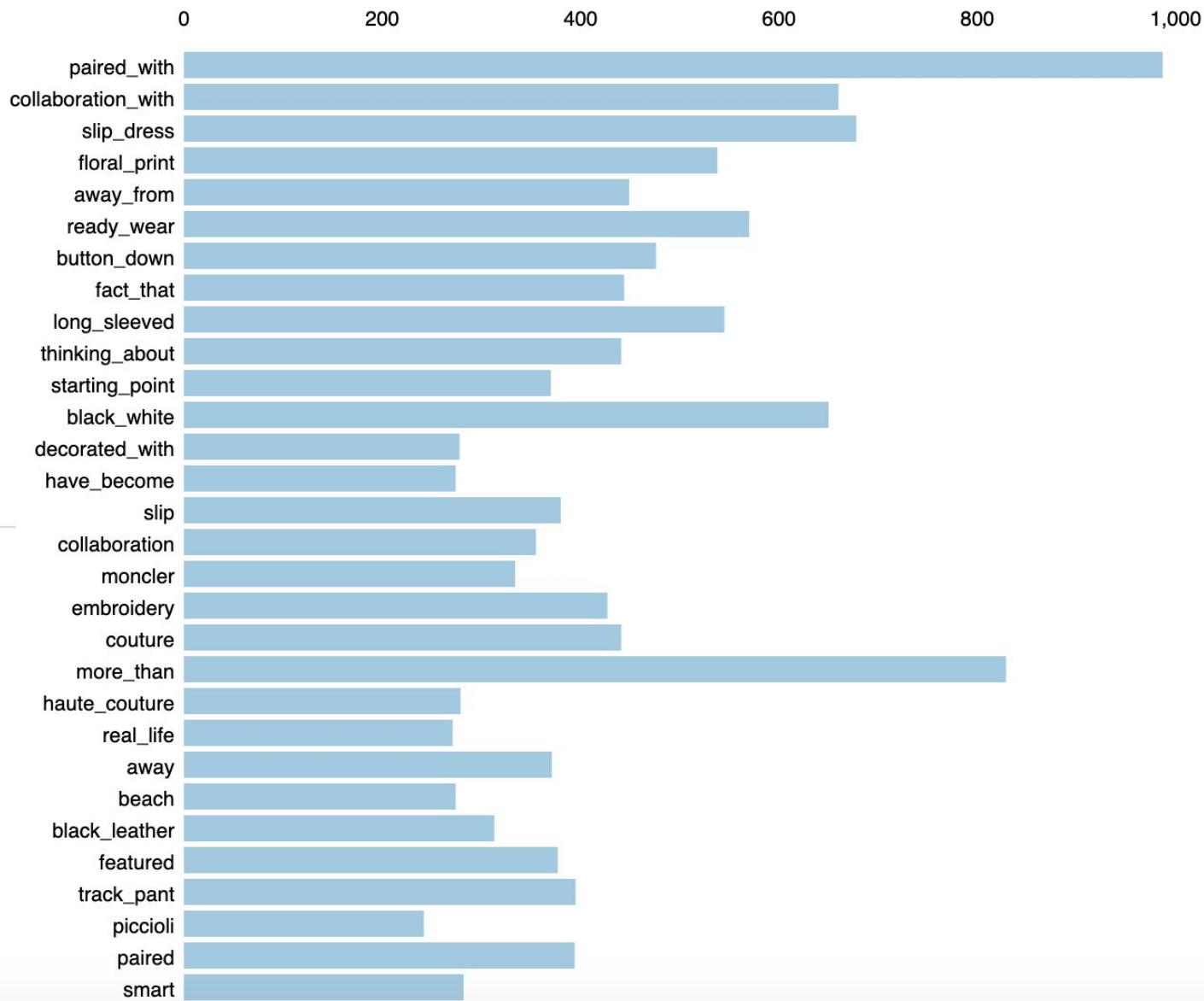
# Choosing number of Topics

- Tried a few more models with various number of topics between 10 and 20
- Arrived at 15 where coherence score actually increased
  - Coherence score was 0.403
- This was chosen as the final model

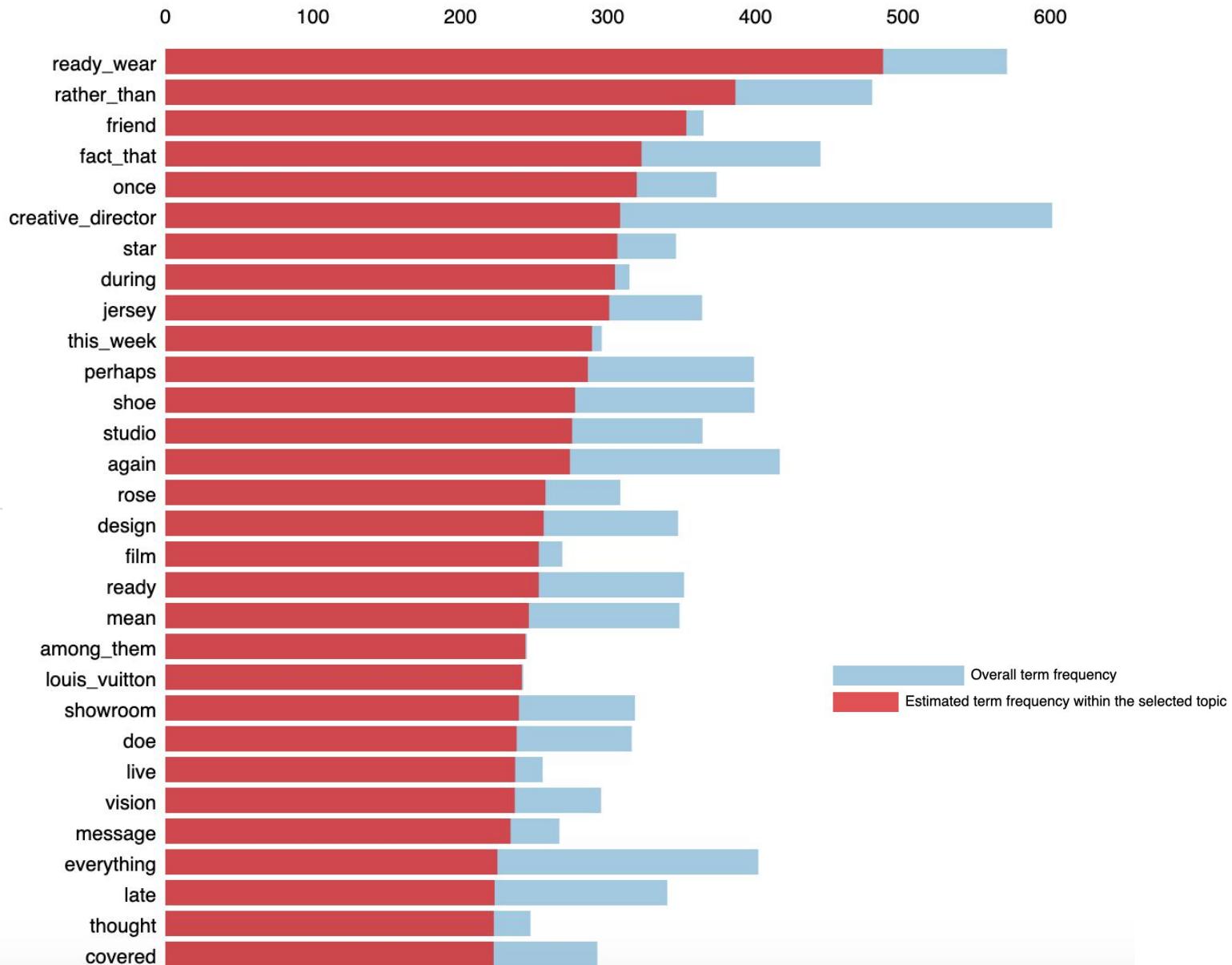
# Topic Relationships



## Top-30 Most Salient Terms<sup>1</sup>



## Top-30 Most Relevant Terms for Topic 1 (30% of tokens)



## Topic 1



*Jersey*



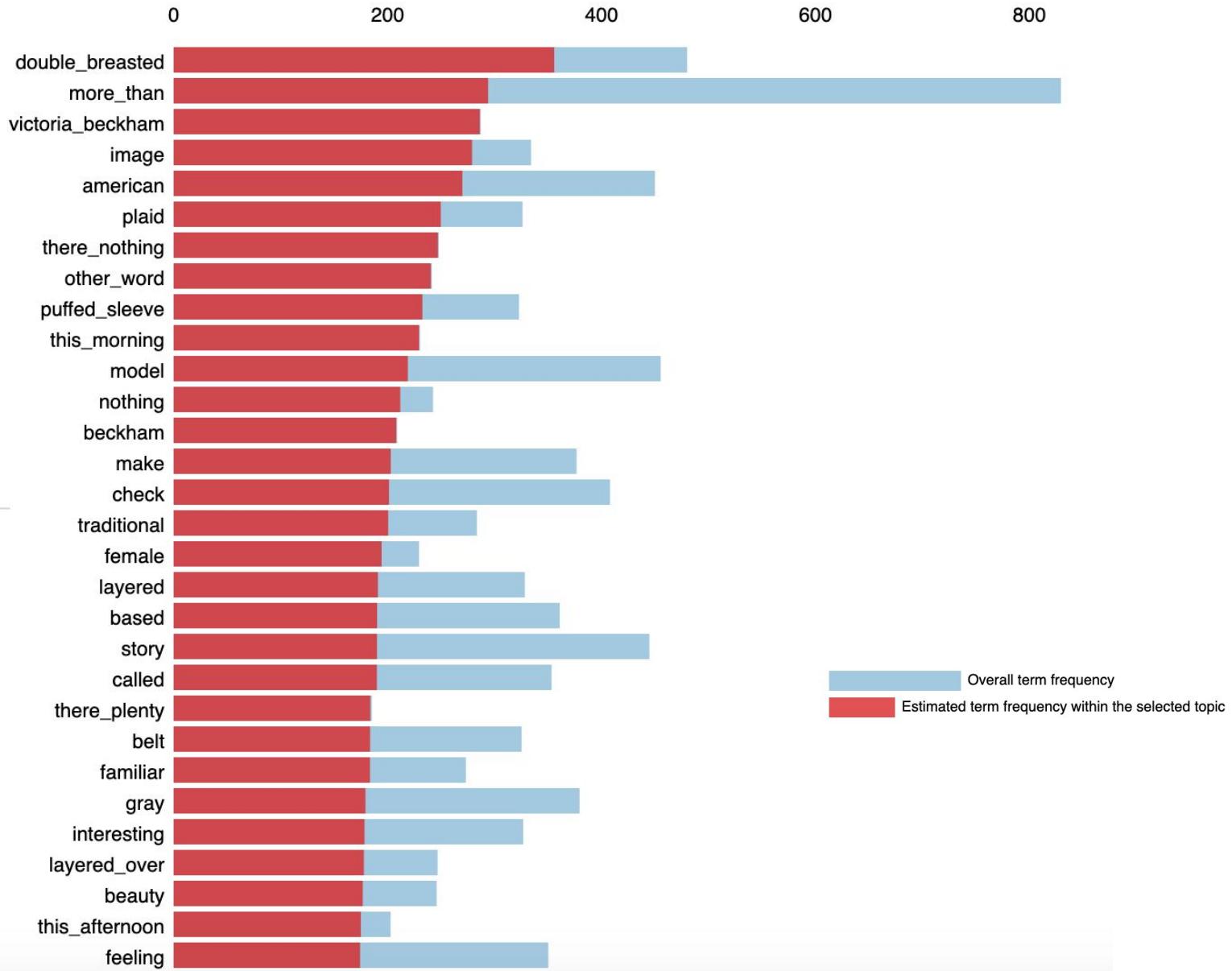
*Or is this Jersey?*



*Louis Vuitton*

Has mostly language terms,  
and proper nouns, not  
specifics

## Top-30 Most Relevant Terms for Topic 2 (18.5% of tokens)



## Topic 2



*Double Breasted*



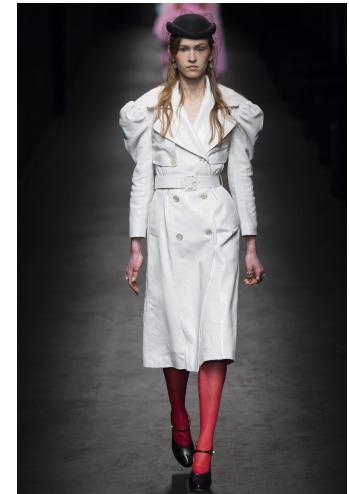
*Plaid & Double  
Breasted*



*Plaid & Puff Sleeve*



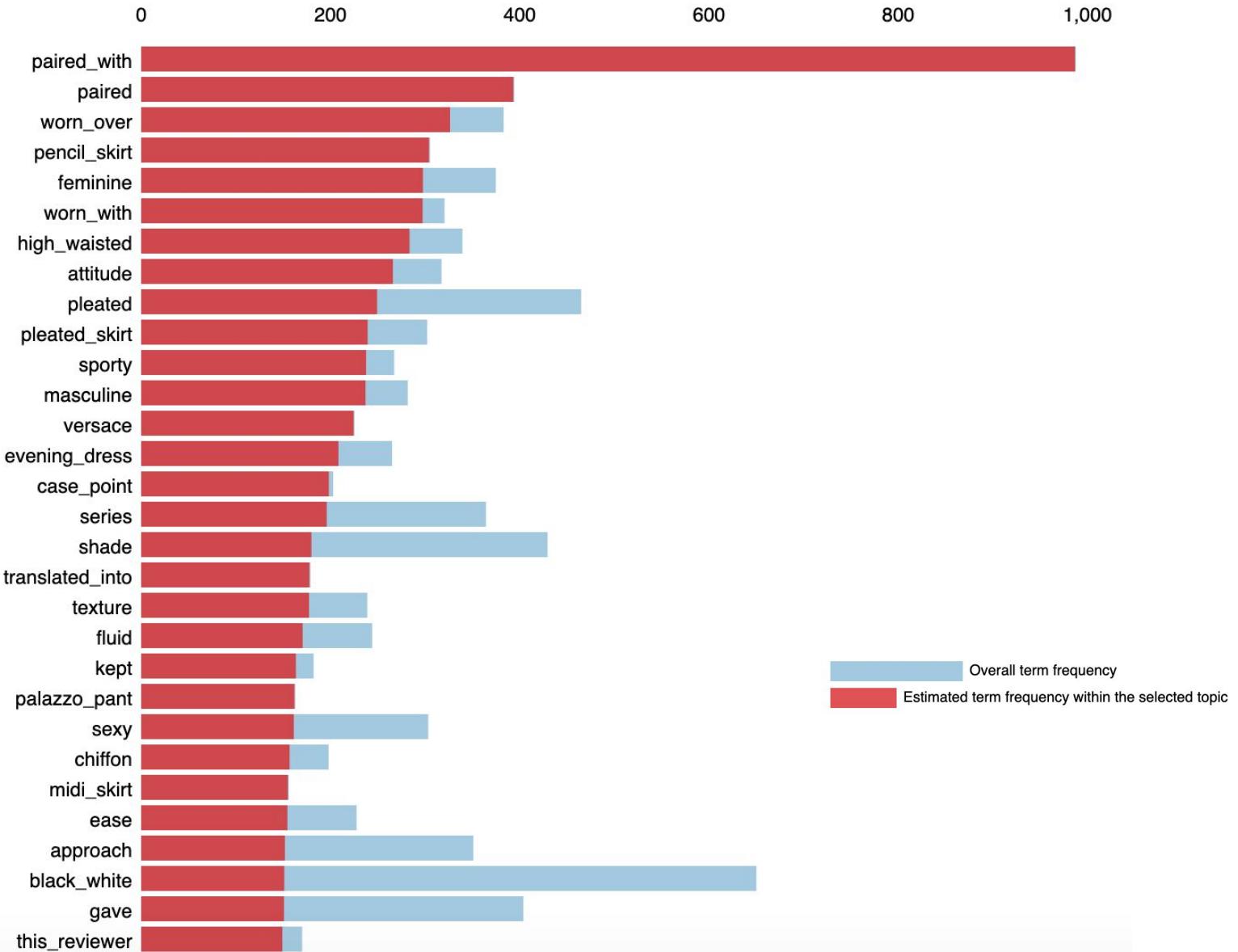
*Puff Sleeve*



*Double Breasted & Puff  
Sleeve*

Many of the trends intersect

### Top-30 Most Relevant Terms for Topic 3 (10.9% of tokens)



## Topic 3



***Pencil Skirt***



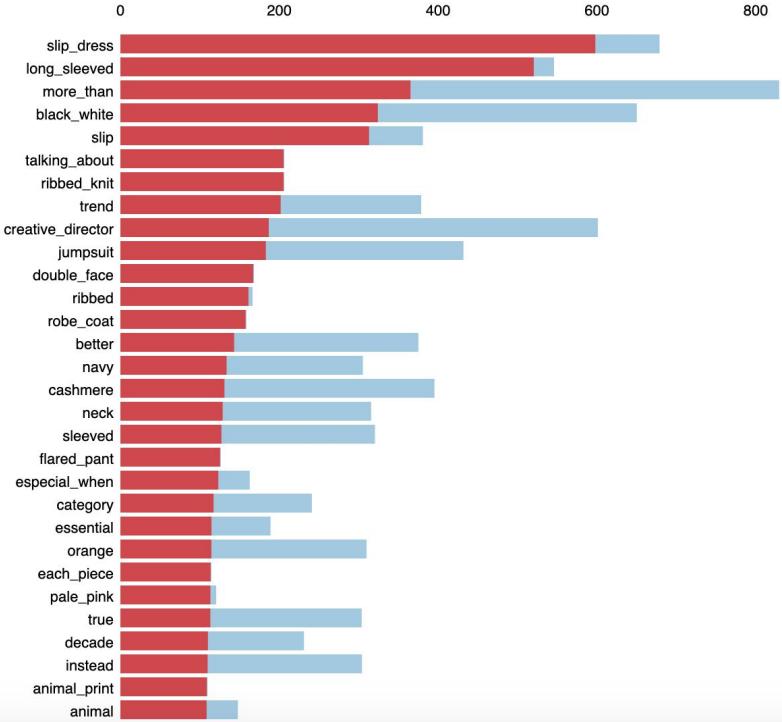
***High Waisted***



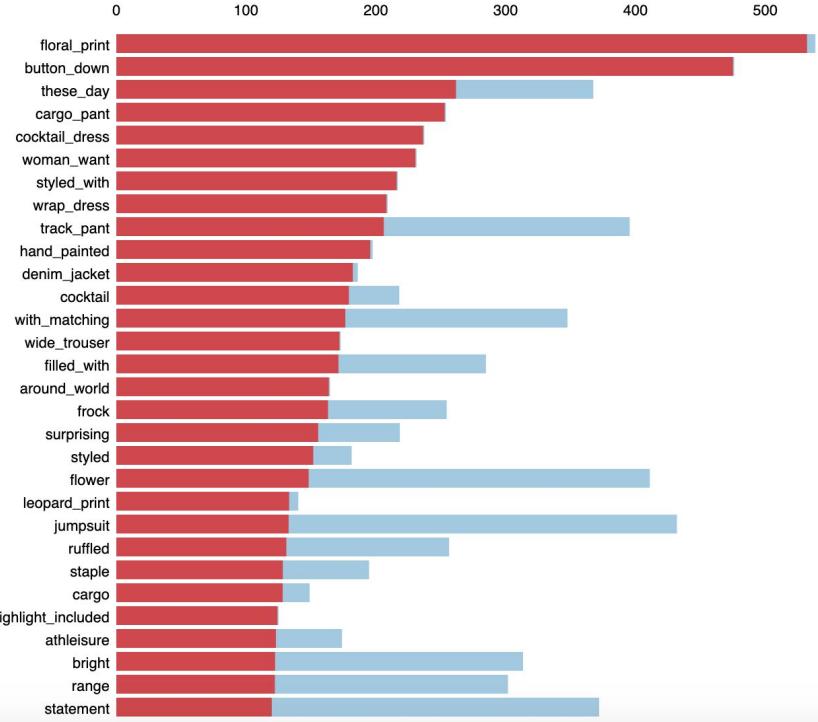
***Pleated Skirt***

Bigrams proved important for trends

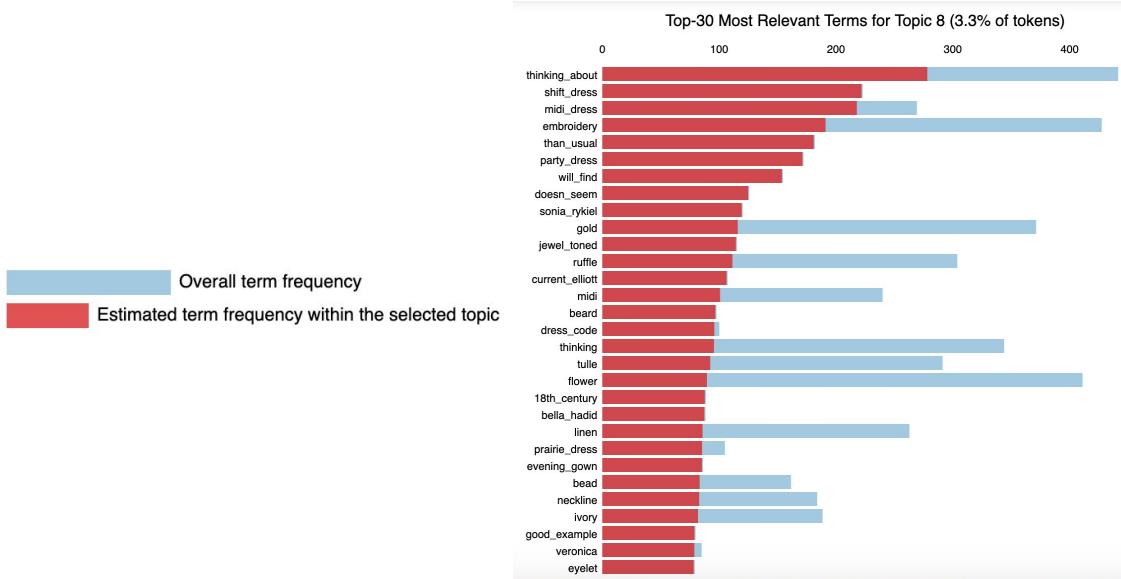
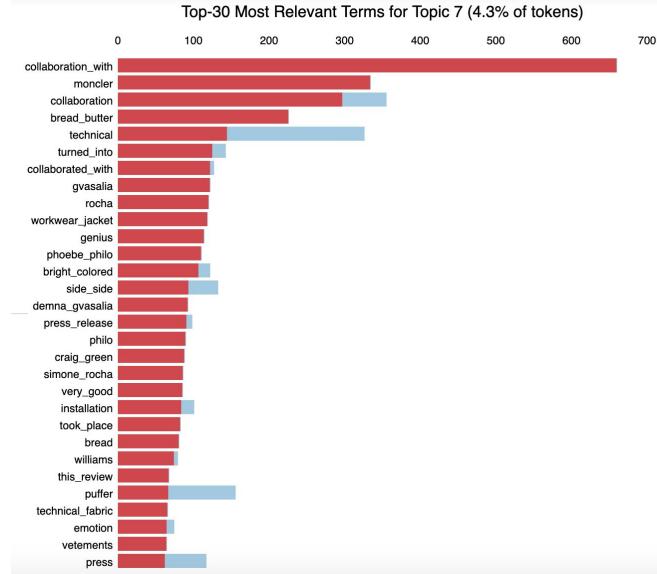
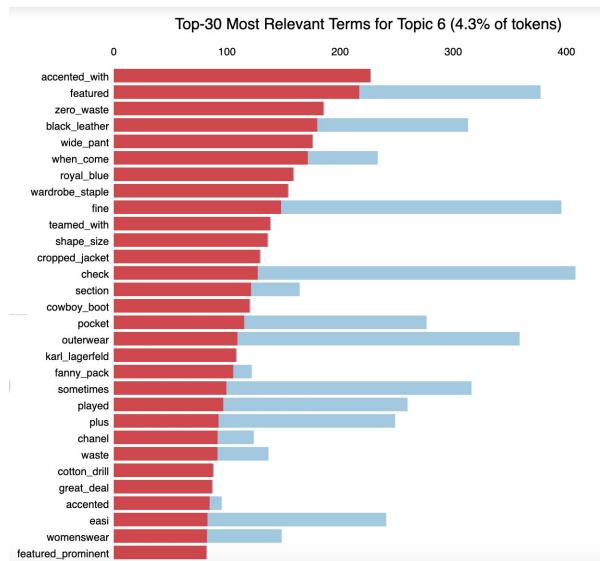
Top-30 Most Relevant Terms for Topic 4 (8.1% of tokens)



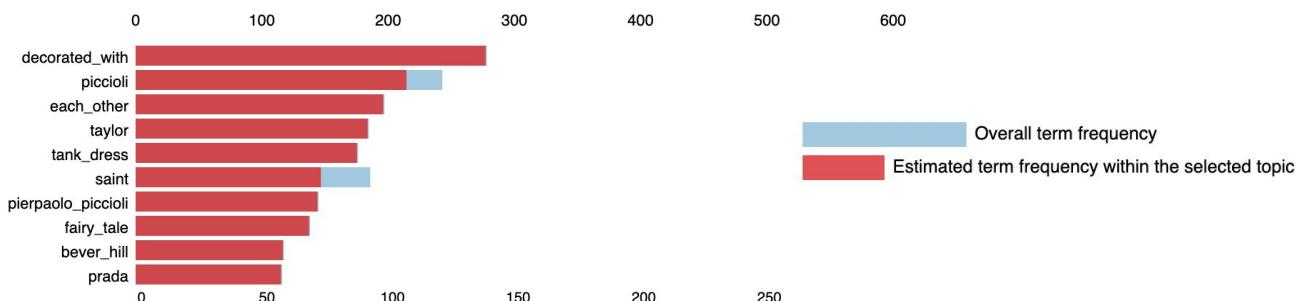
Top-30 Most Relevant Terms for Topic 5 (6.1% of tokens)



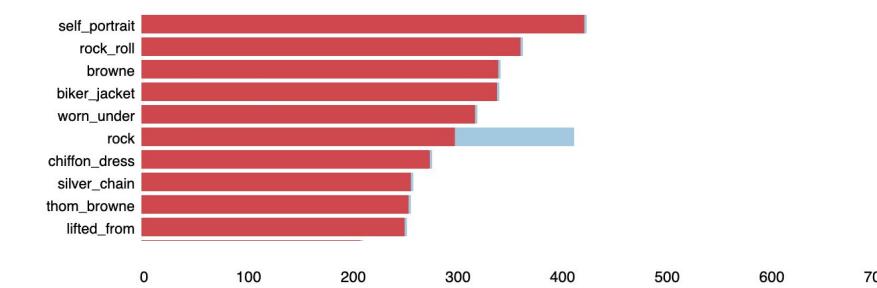
Overall term frequency  
Estimated term frequency within the selected topic



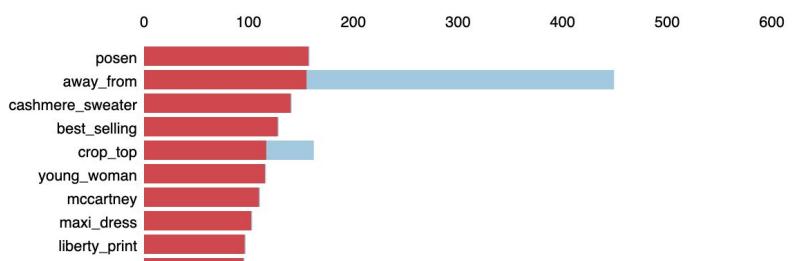
## *Topic 12* (2.1% of tokens)



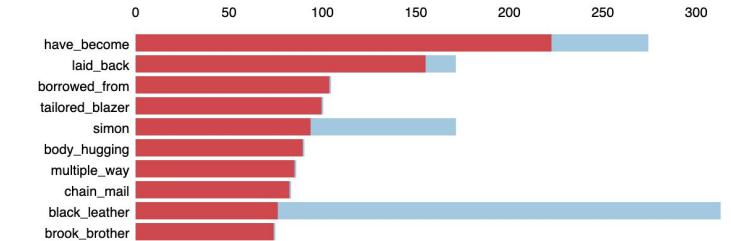
## *Topic 13* (1.91% of tokens)



## *Topic 14* (1.4% of tokens)



## *Topic 15* (1.31% of tokens)



# *Future Work:*

- Analyze text more for additional stop words to eliminate designer names and common words not relevant to goals of this topic modeling
- Collect even more data
- Figure out ways to further improve coherence score.
- Create bar plot visual for topic distribution
- Test these topics on other vogue runway text



*thanks!*

Any questions?