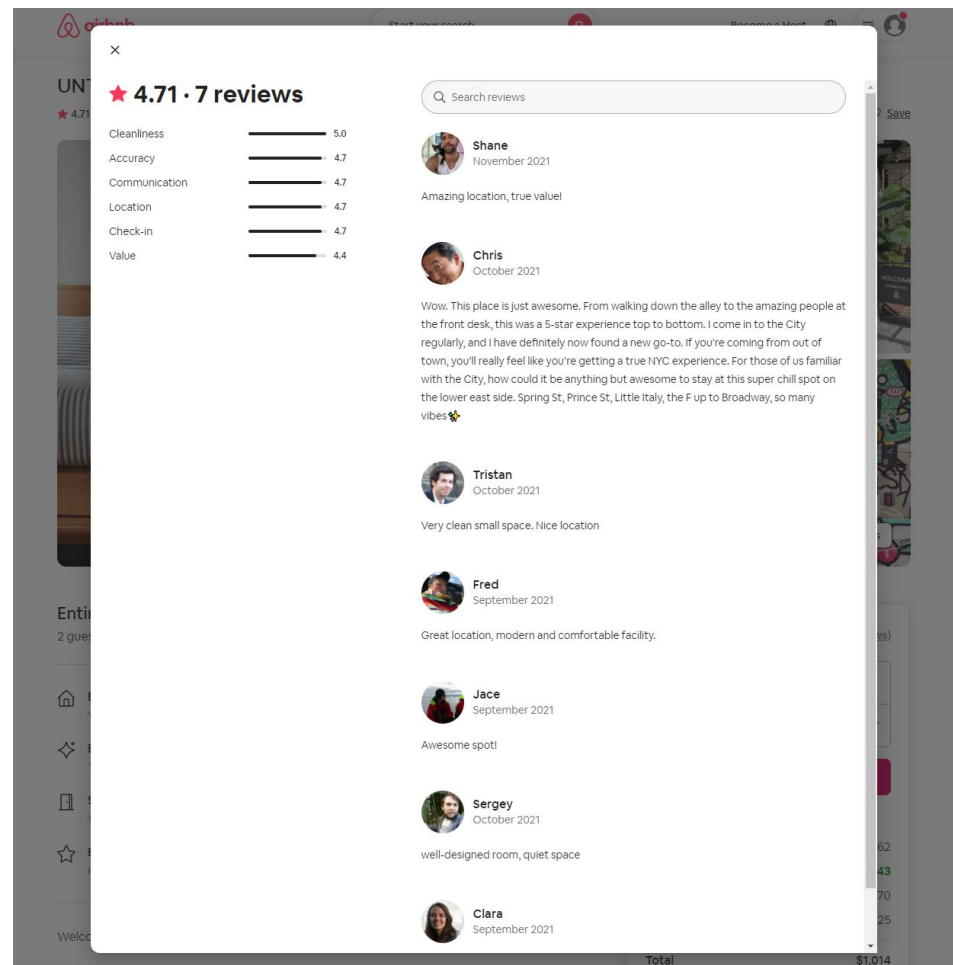




## **Topic Modelling for Reviews, Sources Affecting Guests' Positive and Negative Experiences**

# Why reviews matter?

- Reviews **help guests** choose their travel plans wisely
- Reviews enable hosts to open their homes with confidence and **attract guests.**



# Why reviews matter?

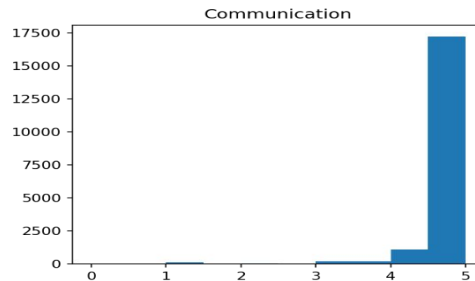
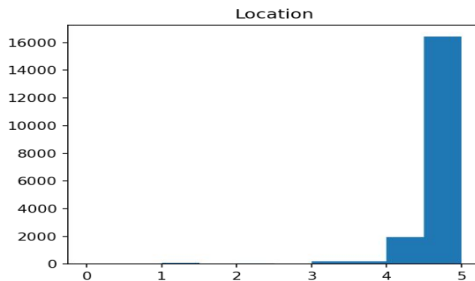
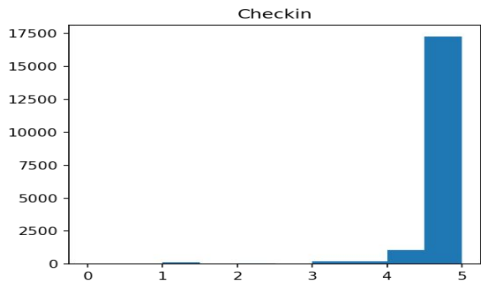
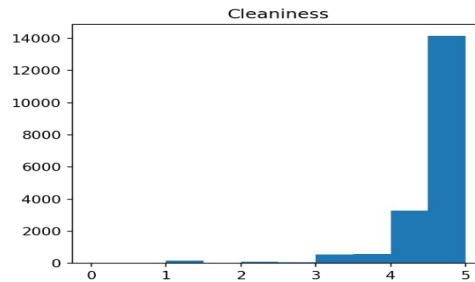
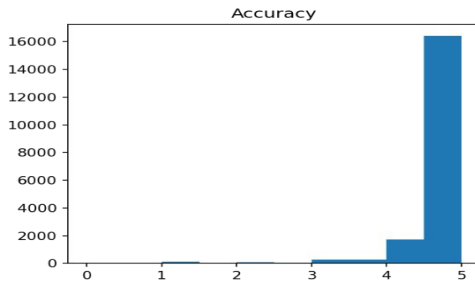
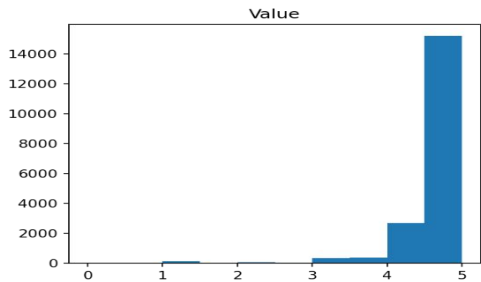
- The first thing to avoid when booking an Airbnb:

## Little or No Reviews



# What about ratings?

- Guests can provide star ratings with 1–5 stars.
- All review scores are highly positive scores.
- There is no scores less than 4.5 out of 5.



# Methodology

## Data

- NYC Listings Dataset, InsideAirbnb
- More than 80K Reviews

## Data Cleaning & EDA

- Remove numbers, capital letters and punctuations
- Eliminate non-English reviews
- Lemmatize
- Tools: Pandas, Numpy, langdetect, NLTK, Matplotlib, Seaborn

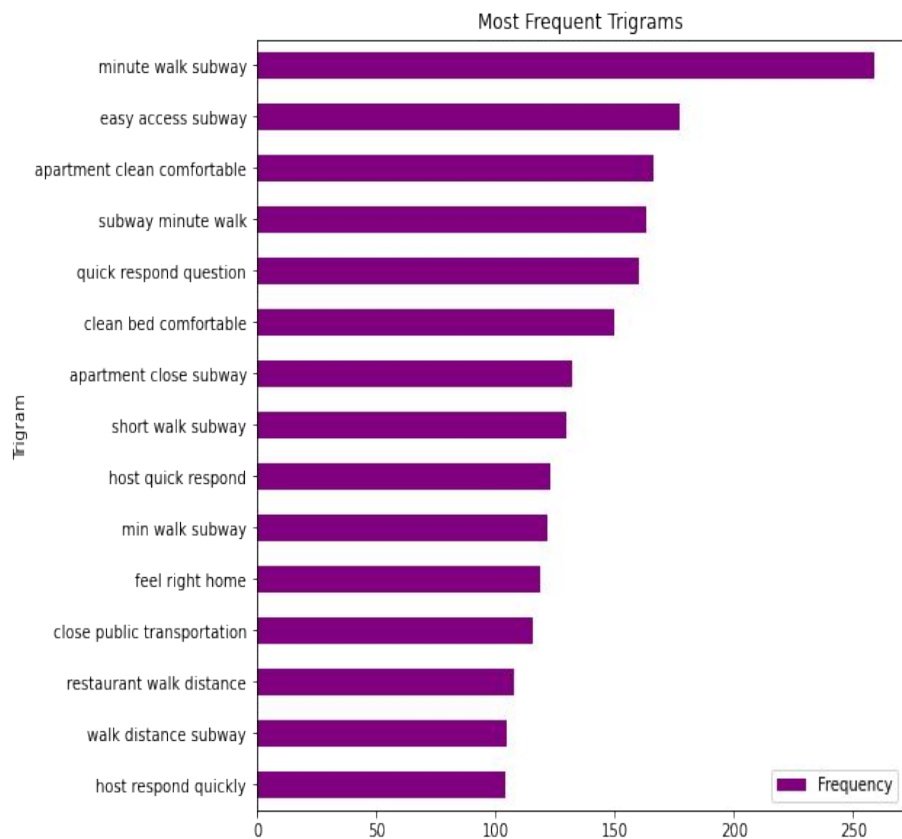
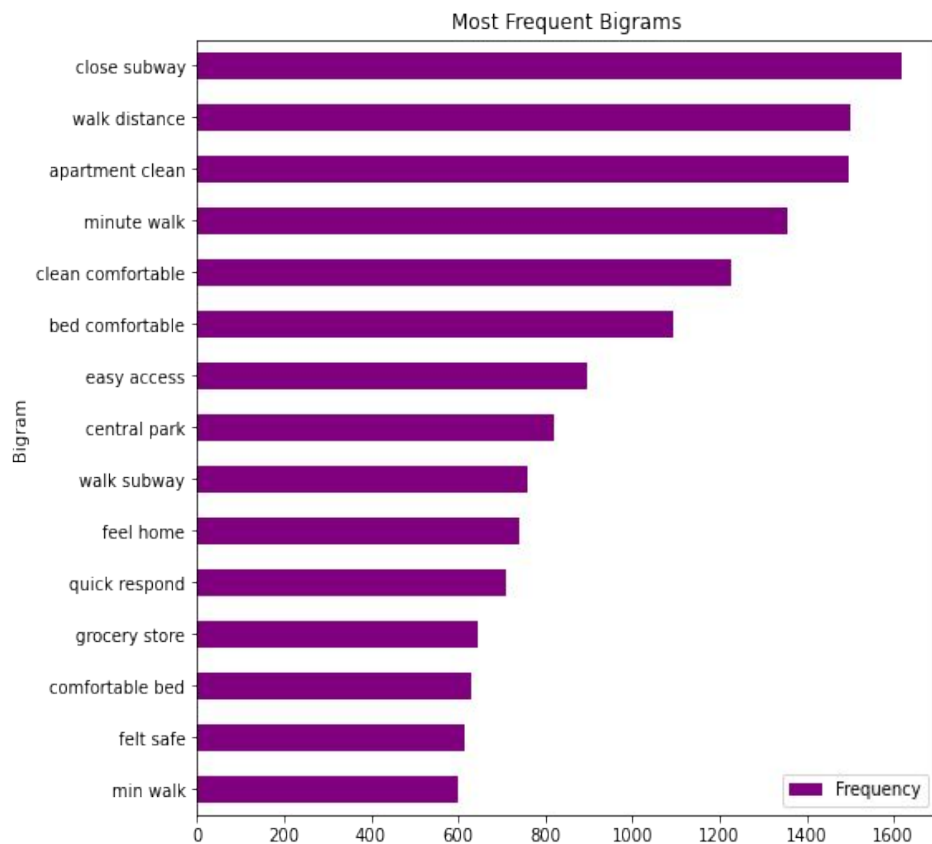
## Topic Modeling

- TF-IDF Vectorizer
- CountVectorizer
- Models :
  - NMF
  - SVD
  - LDA
  - CorEx
- Tools : sklearn, pyLDAvis

## Final Model

- Count Vectorizer (ngram\_range=(1,2), max\_df = 0.7, min\_df = 10)
- LDA
- VaderSentimentAnalysis
- Tools : sklearn, pyLDAvis

# Reviews: Most common words in reviews



# Topic Modelling

- Vectorizer :
  - CountVectorizer
- Topic Modeler :
  - LDA
- Number of Topic 15

## Topics

---



Rental interior



Kitchen



Overall Airbnb Experience



Neighborhood- accessibility to transportation



Home-like comfort/experience



Neighborhood/ accessibility to dining, social attractions



Cleanliness



Host-hospitality



Location- safe/family friendly



Bed/Bathroom



Overall trip experience



Host-responsiveness



Convenience (check in/out, comfort, hotel like)



Comfort/Value

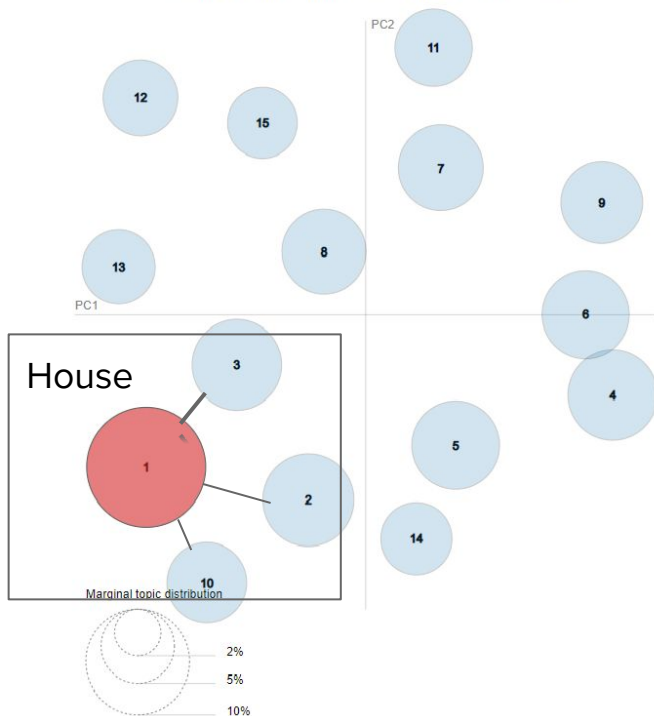


Listing Accuracy

# Topic Visualization

Selected Topic:

Intertopic Distance Map (via multidimensional scaling)

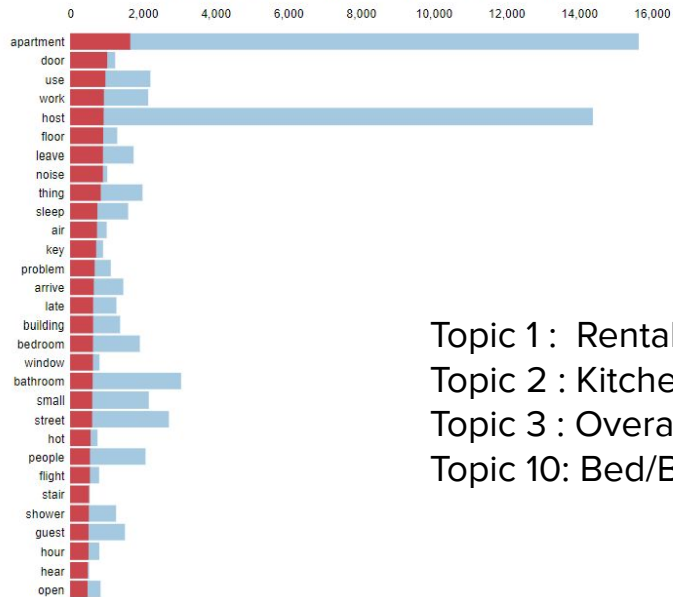


Slide to adjust relevance metric:<sup>(2)</sup>

$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Top-30 Most Relevant Terms for Topic 1 (13.2% of tokens)



Overall term frequency

Estimated term frequency within the selected topic

1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

Topic 1 : Rental interior

Topic 2 : Kitchen

Topic 3 : Overall experience

Topic 10: Bed/Bathroom



# Topic Visualization

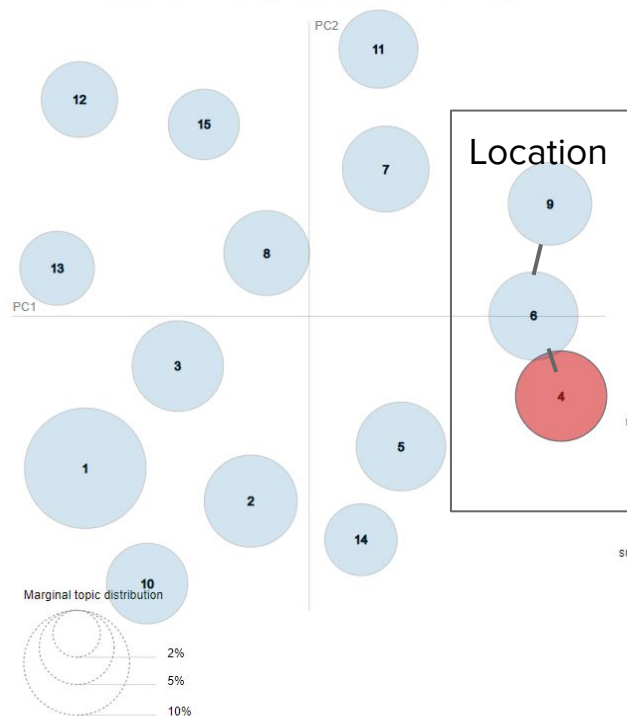
Selected Topic:  Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>

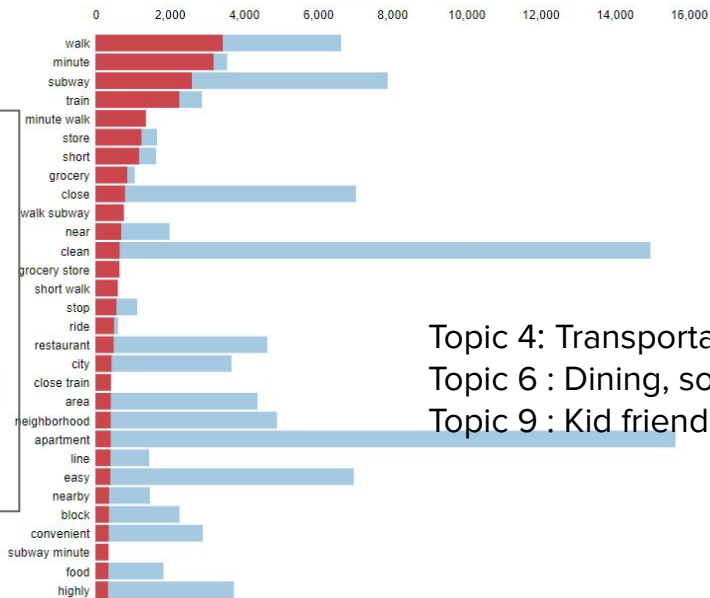
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 4 (7.4% of tokens)



Topic 4: Transportation

Topic 6 : Dining, social attractions

Topic 9 : Kid friendly, safe

Overall term frequency  
Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

# Topic Visualization

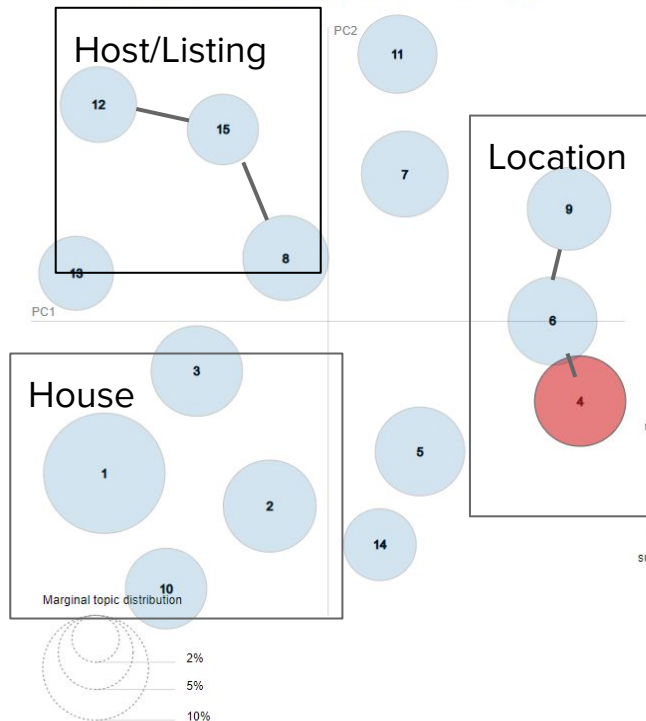
Selected Topic:  Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>

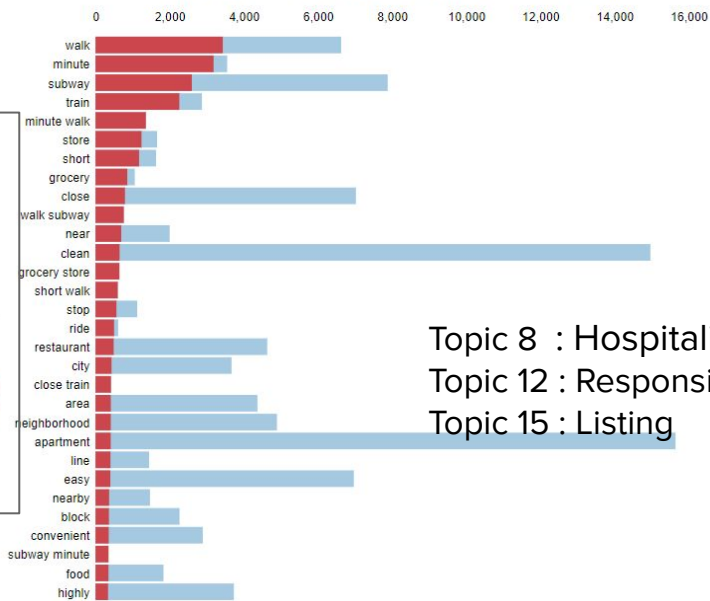
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 4 (7.4% of tokens)

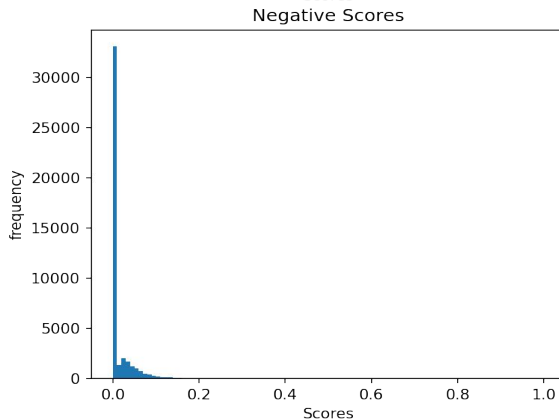
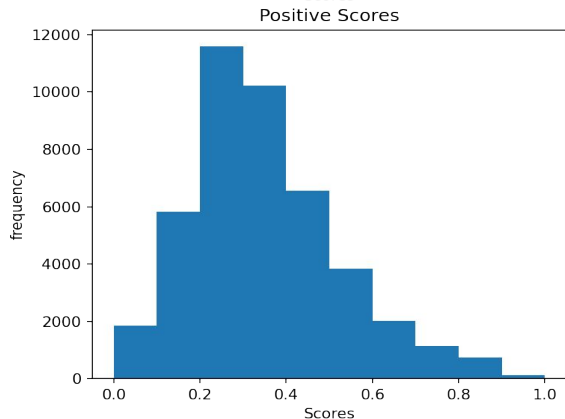
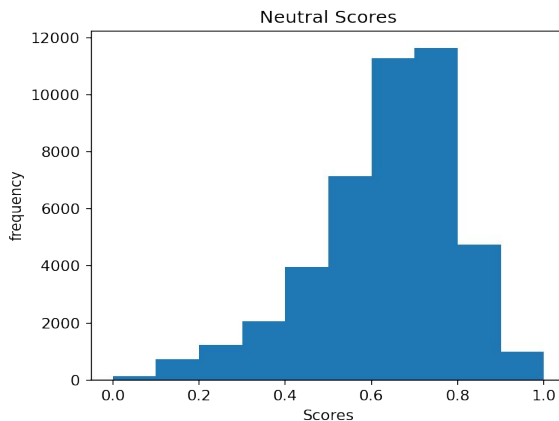
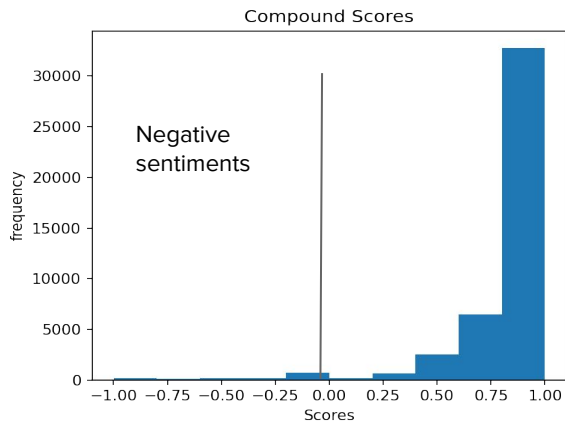


Topic 8 : Hospitality  
Topic 12 : Responsiveness  
Topic 15 : Listing

Overall term frequency  
Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)  
2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

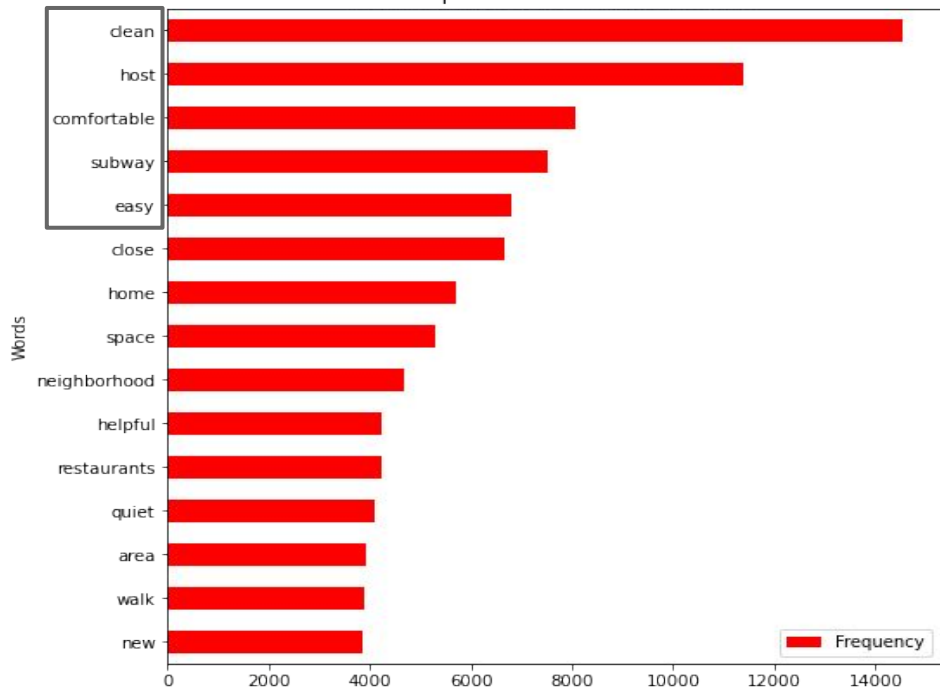
# Reviews: Sentiment Analysis/ VaderSentiment



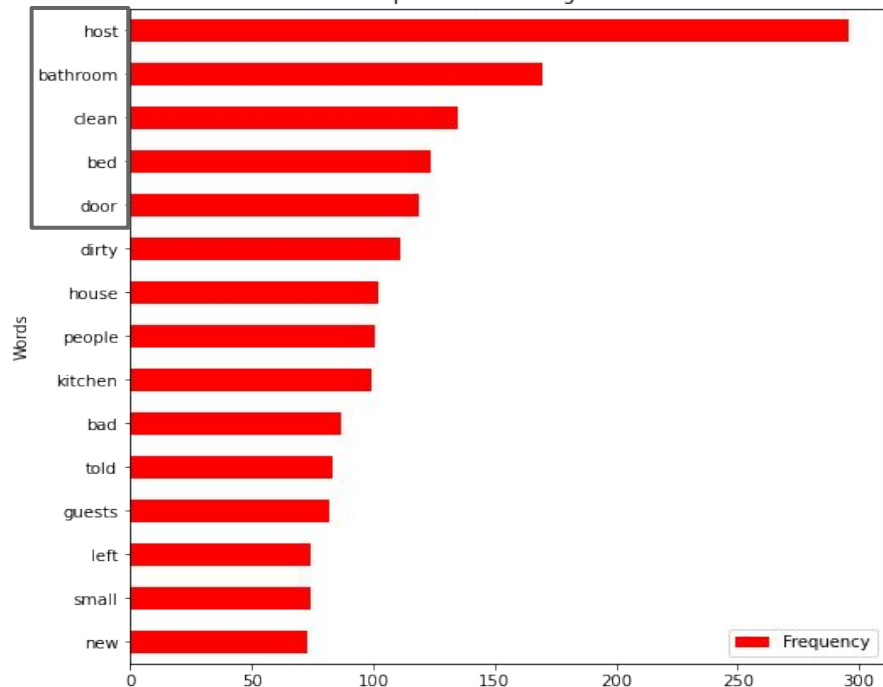
- Most guest had pleasant experiences during their Airbnb stays.

# Positive vs Negative Reviews

Most Frequent Words - Positive Reviews



Most Frequent Words - Negative Reviews



Host interaction, Cleanliness → Overall experience

# Insights for Hosts

## Common topics in reviews:

- The condition of the house
  - Kitchen, bedrooms, and bathrooms
- Cleanliness
- Location
- Host responsiveness and hospitality

## Next Steps

- Cluster analysis to group topics
- Topic modeling for different geographic districts, different type of listings
- Design a recommendation system where common topics are a filtering option
- Explore effect of reviews to occupancy rates

# Questions



*Thank you!*

# Appendix



## Unigram





# Review Topics

- 1 Apartment interior issues
- 2 Kitchen Experience
- 3 Airbnb Experience
- 4 Neighborhood- accessibility to transportation
- 5 Home-like comfort/experience
- 6 Neighborhood/ accessibility to social attractions
- 7 Cleanliness
- 8 Host-hospitality
- 9 Location- safety/family friendly
- 10 Bed/Bathroom
- 11 Overall trip experience
- 12 Host-responsiveness
- 13 Convenience (check in/out, comfort, hotel like)
- 14 Comfort/Value
- 15 Listing Accuracy

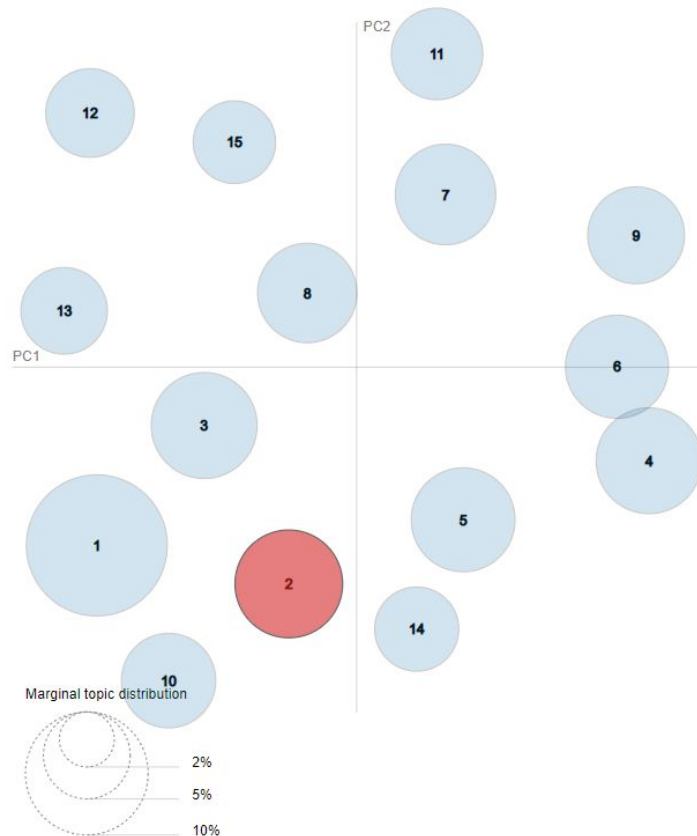
Selected Topic:

Slide to adjust relevance metric:<sup>(2)</sup>

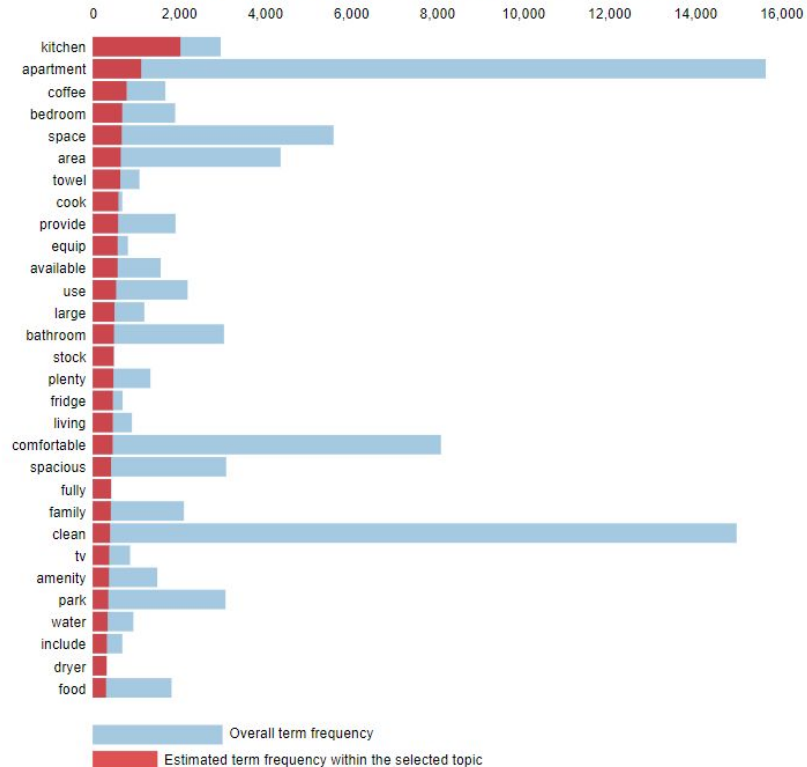
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 2 (7.7% of tokens)



1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Slevert & Shirley (2014)

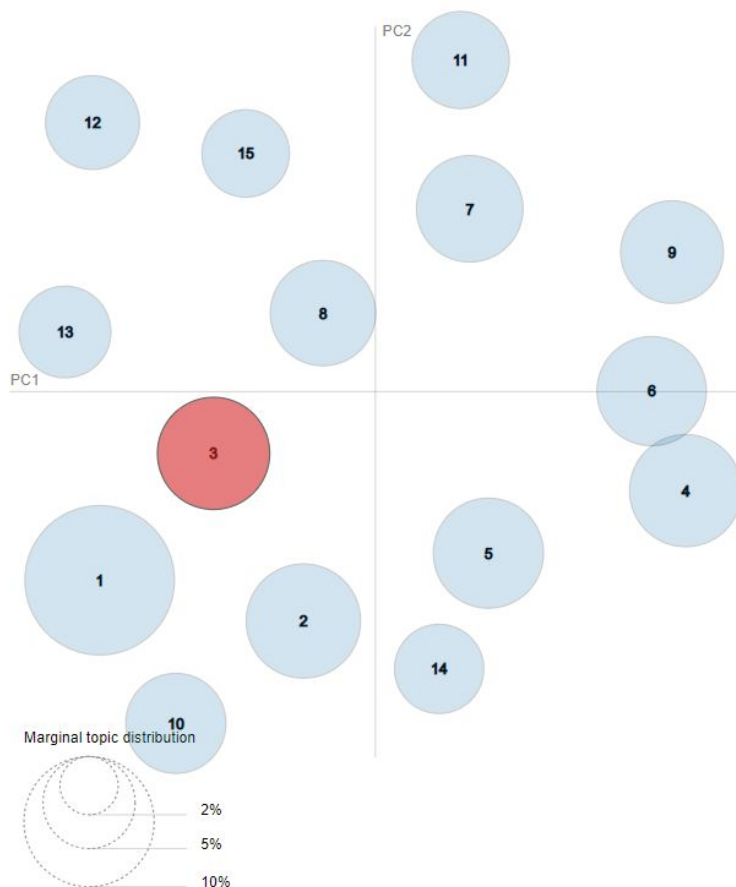
Selected Topic:

Slide to adjust relevance metric:<sup>(2)</sup>

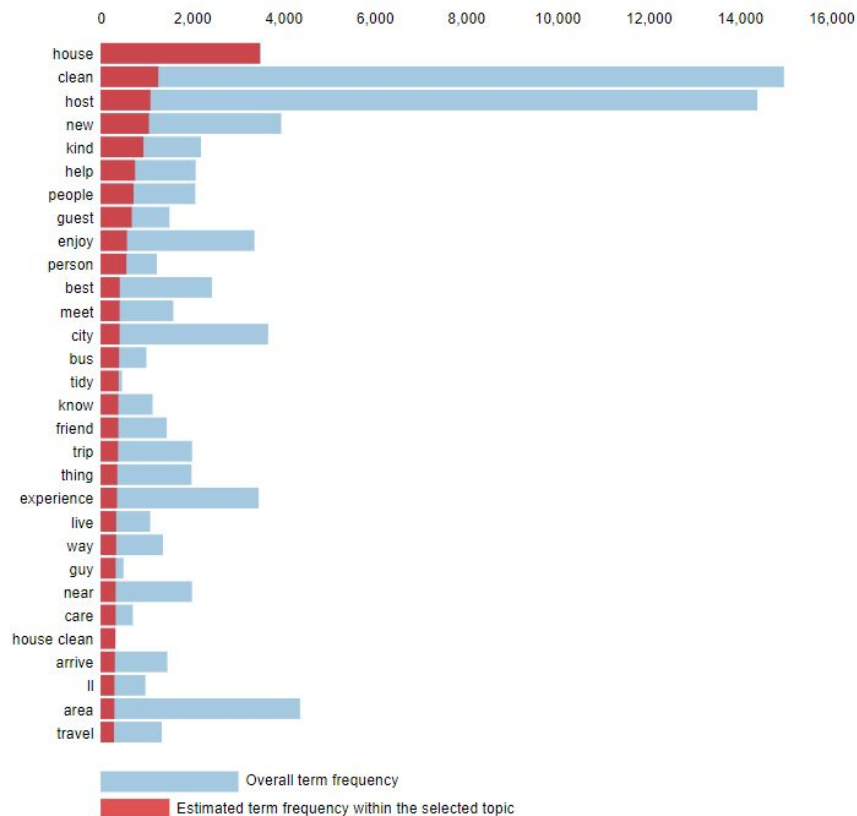
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 3 (7.4% of tokens)



1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

Selected Topic:  Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>

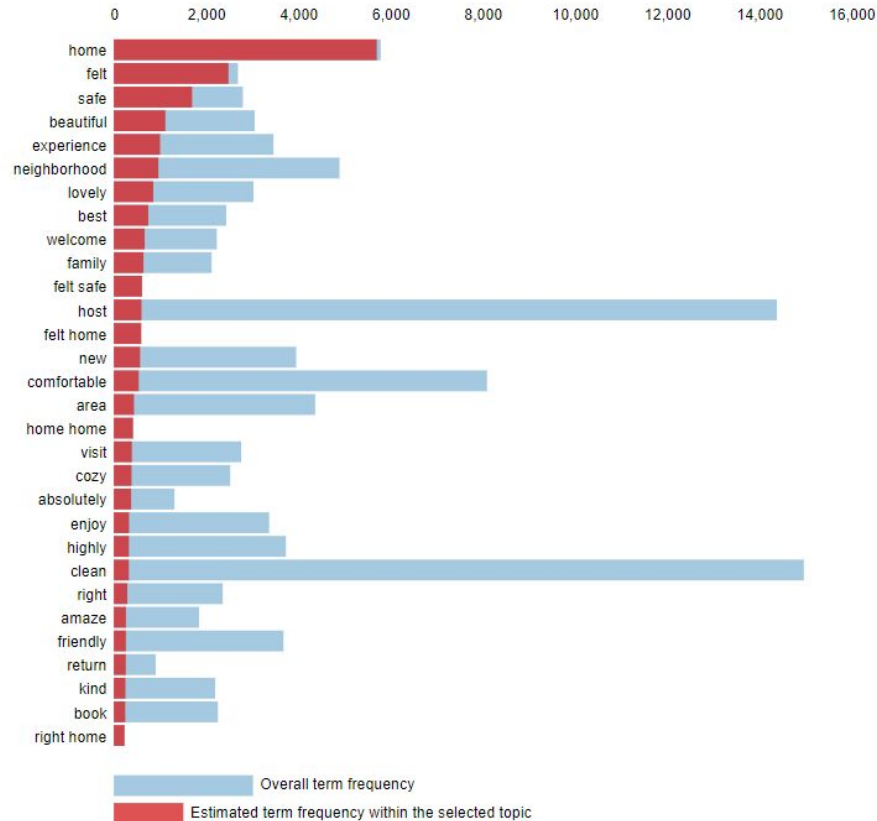
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 5 (7.1% of tokens)



1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)
2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

Selected Topic:

Slide to adjust relevance metric:<sup>(2)</sup>

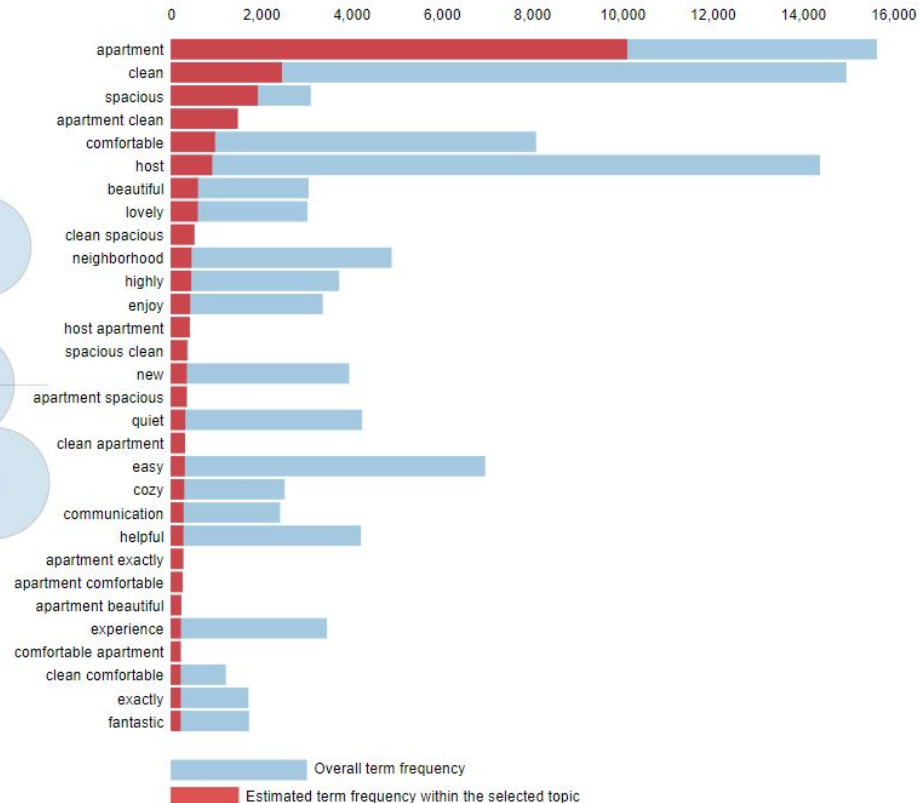
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



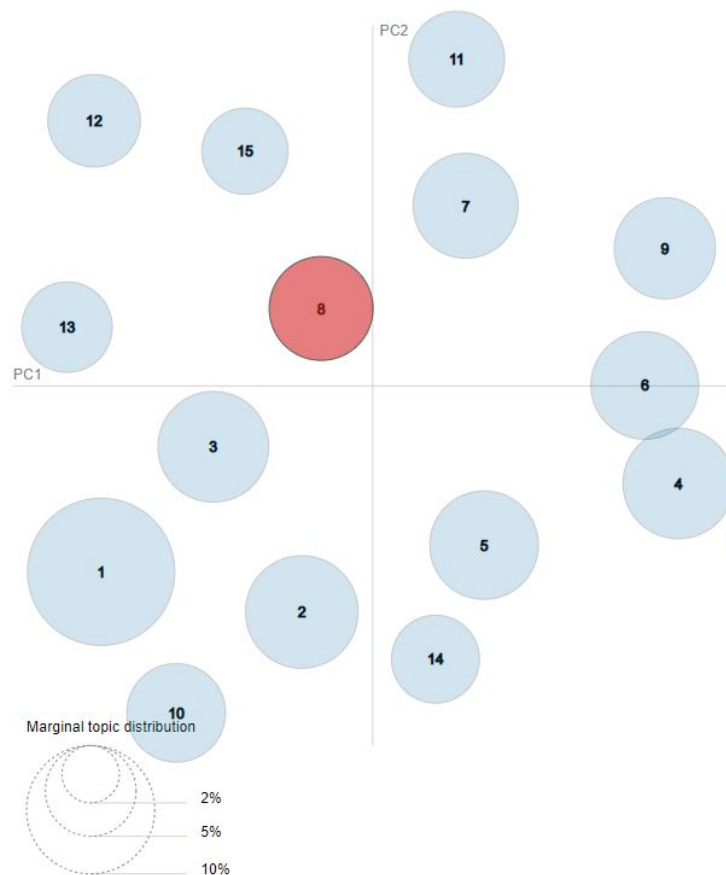
Top-30 Most Relevant Terms for Topic 7 (6.7% of tokens)



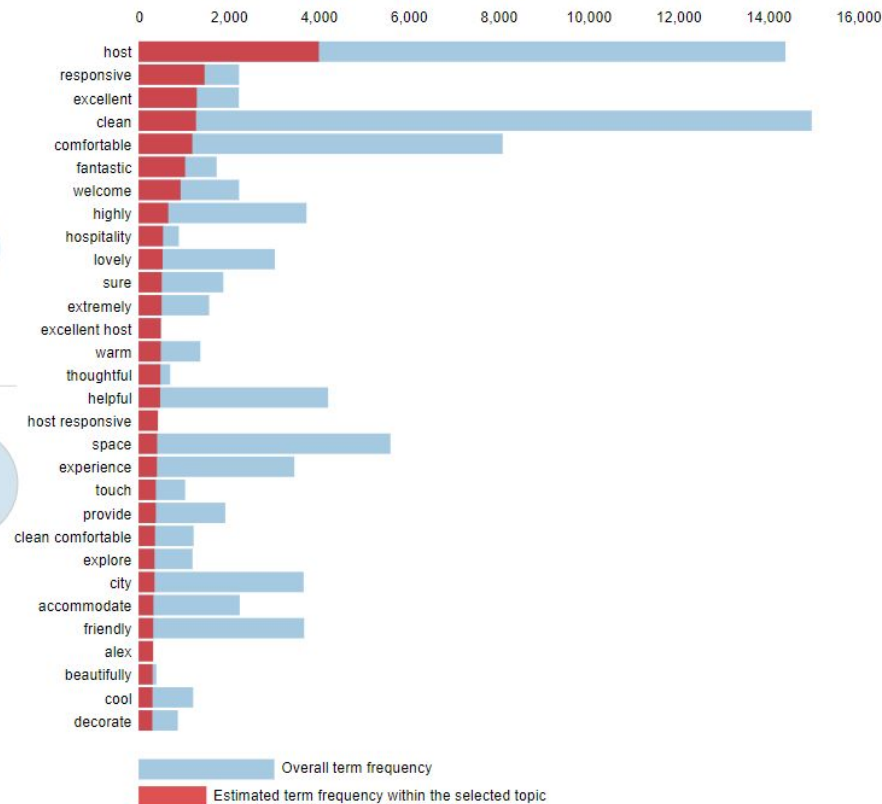
1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) =  $\lambda$  \* p(w | t) + (1 -  $\lambda$ ) \* p(w | t)/p(w); see Sievert & Shirley (2014)

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 8 (6.5% of tokens)



1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)
2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

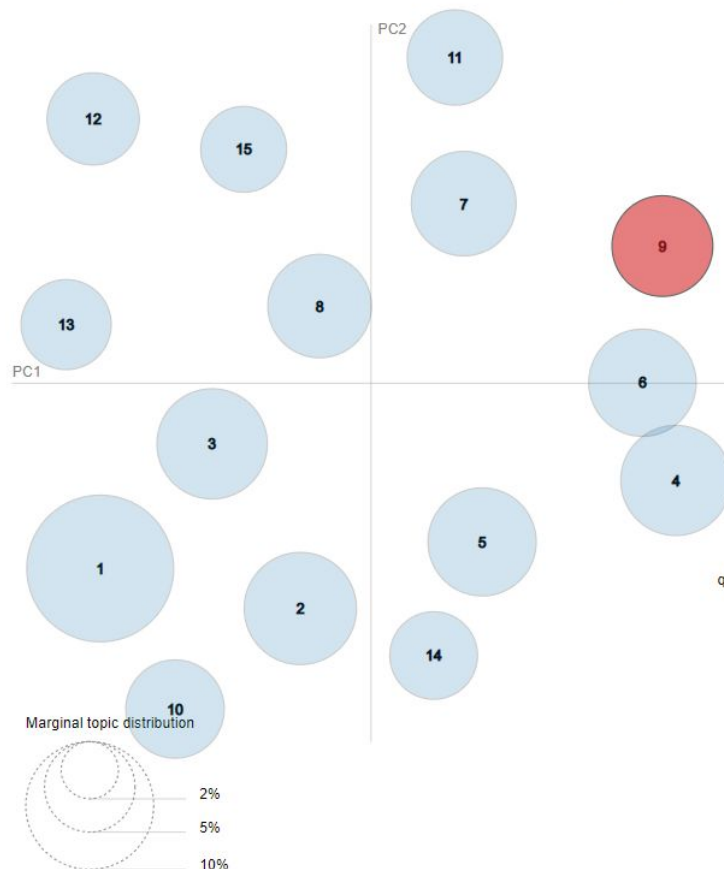
Selected Topic:

Slide to adjust relevance metric:(2)

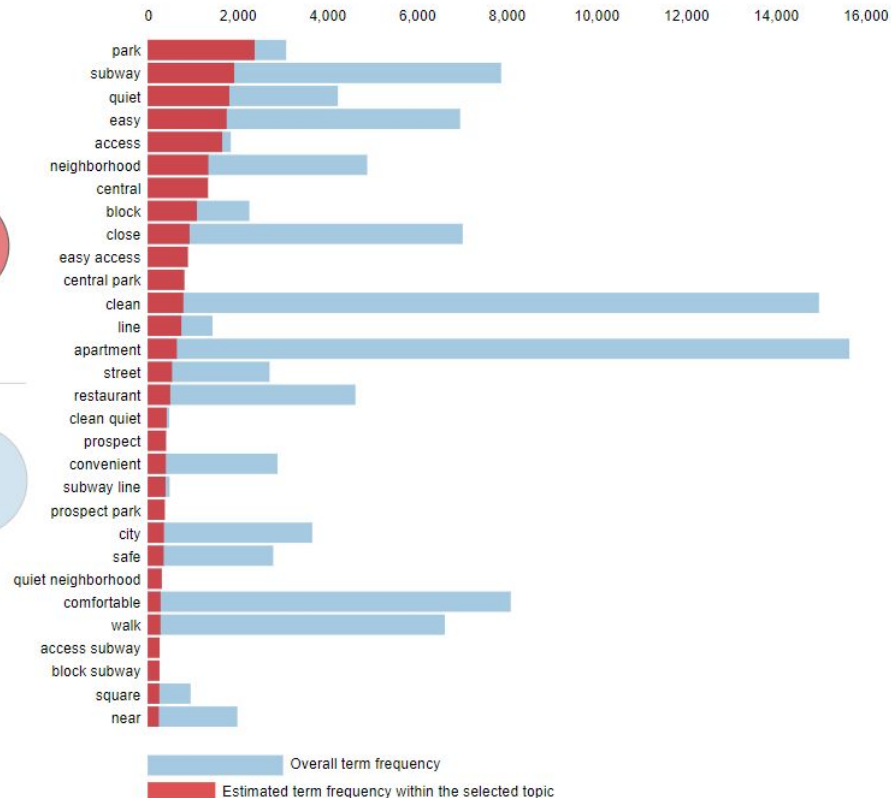
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 9 (6.2% of tokens)

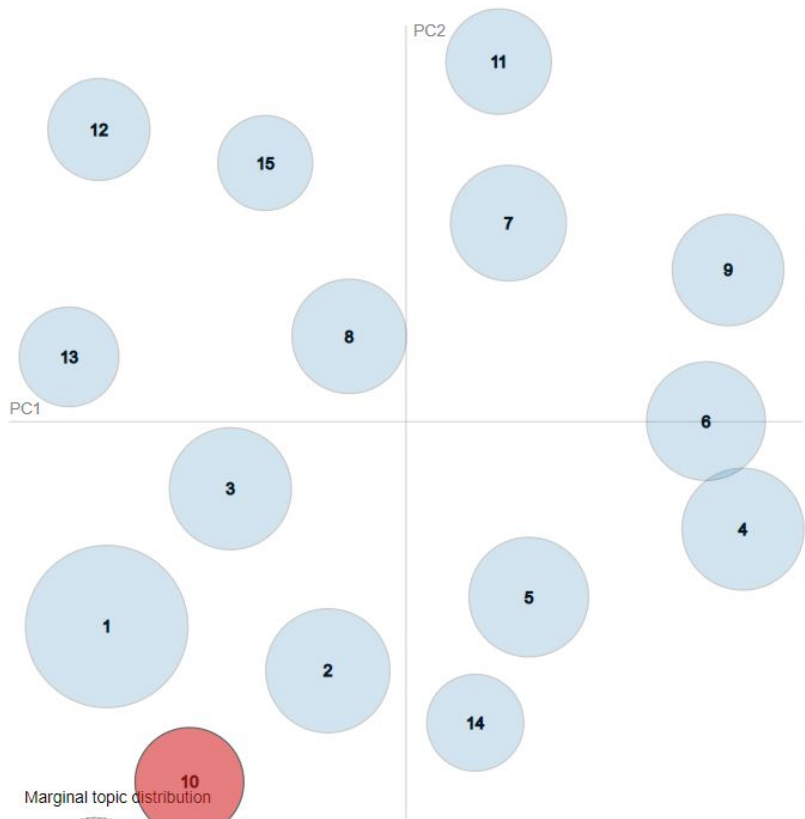


1.  $saliency(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

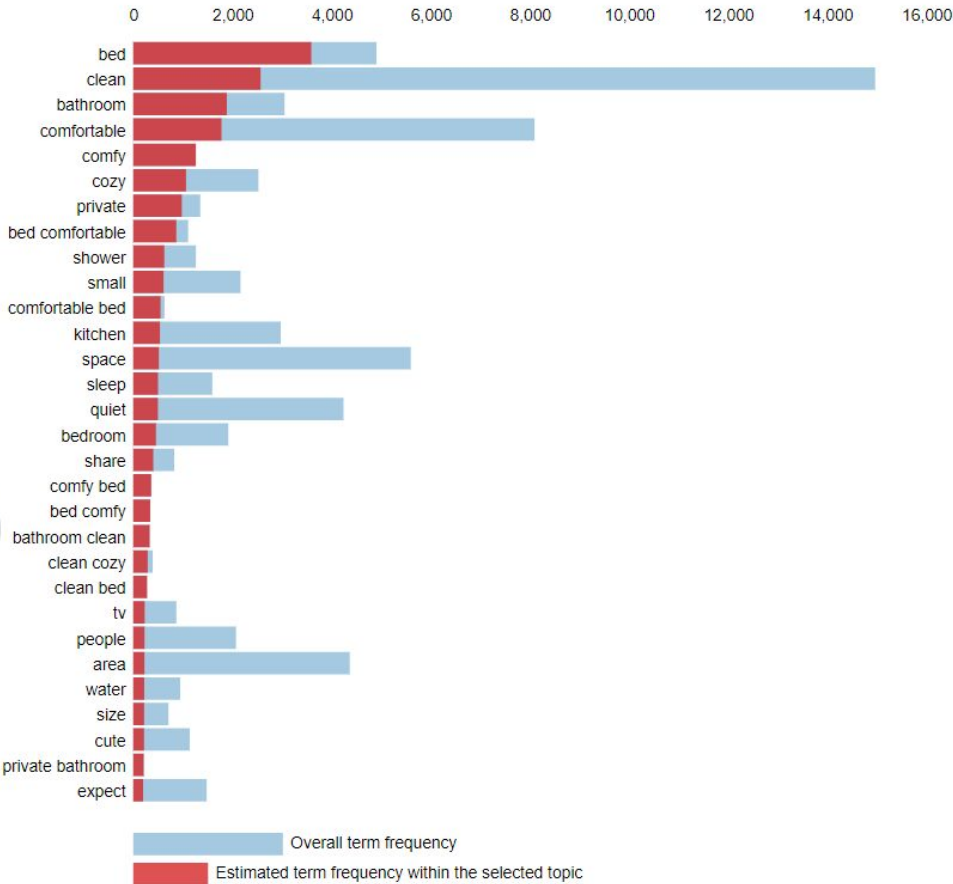
2.  $relevance(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)



Intertopic Distance Map (via multidimensional scaling)

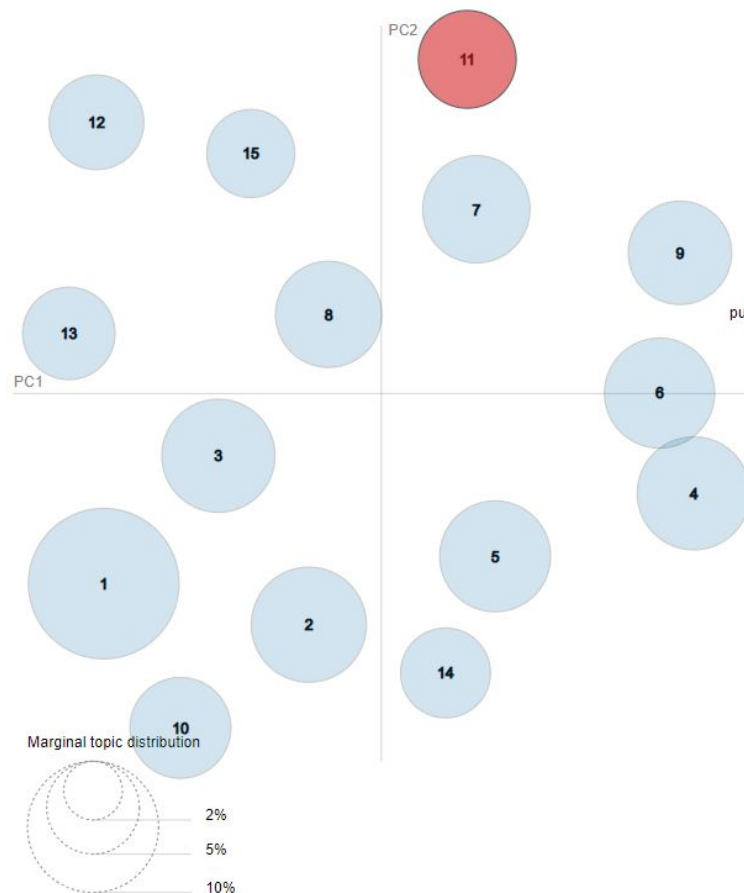


Top-30 Most Relevant Terms for Topic 10 (5.9% of tokens)

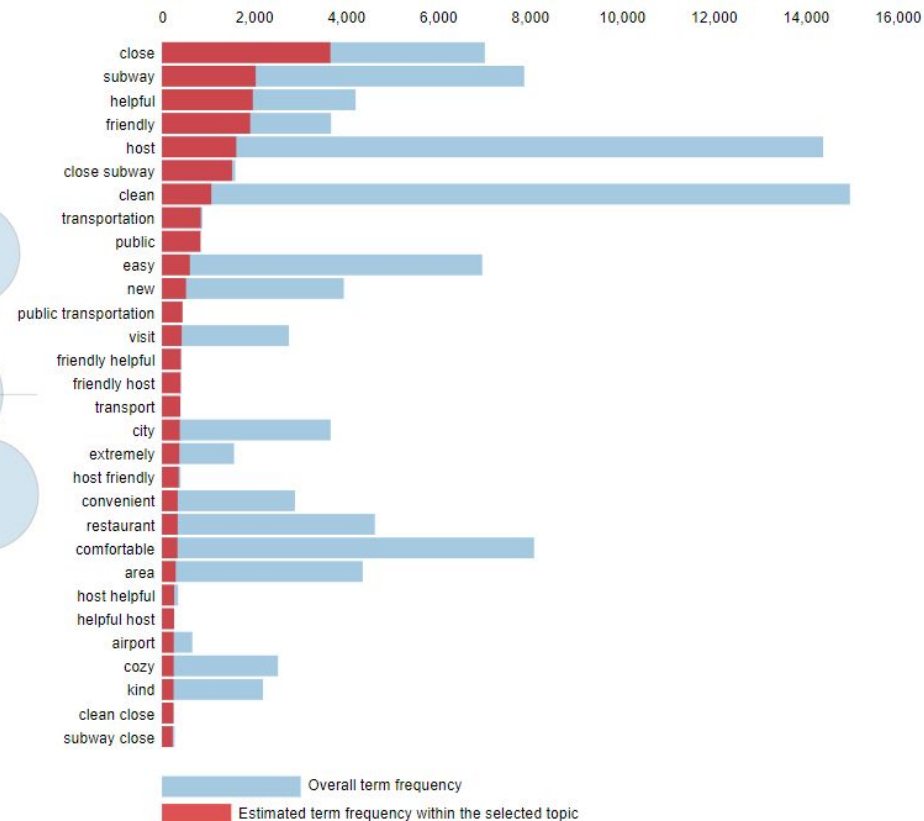




Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 11 (5.5% of tokens)



1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

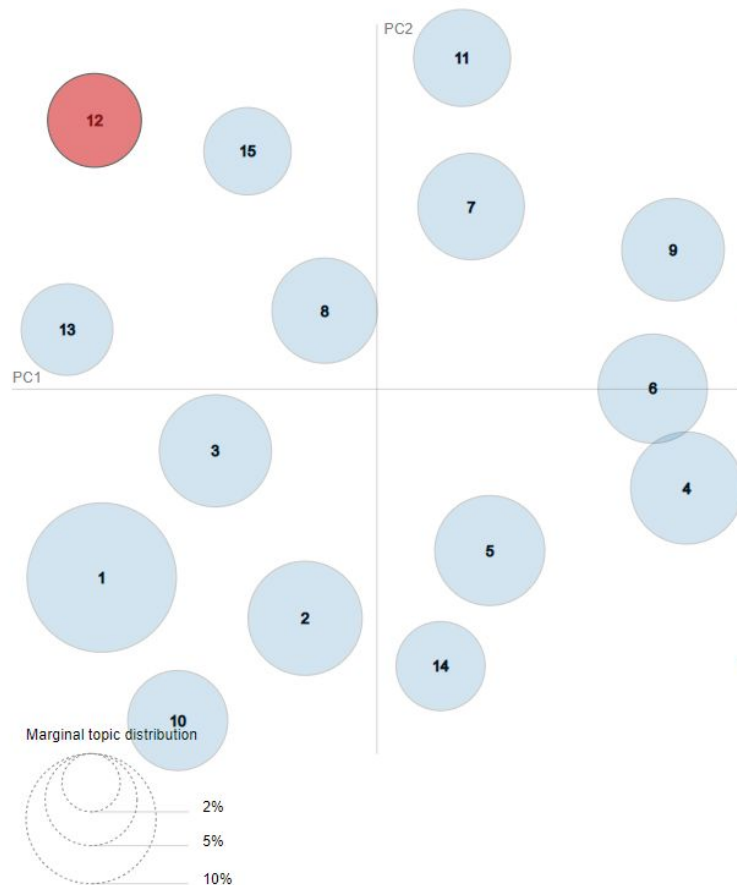
Selected Topic:

Slide to adjust relevance metric:<sup>(2)</sup>

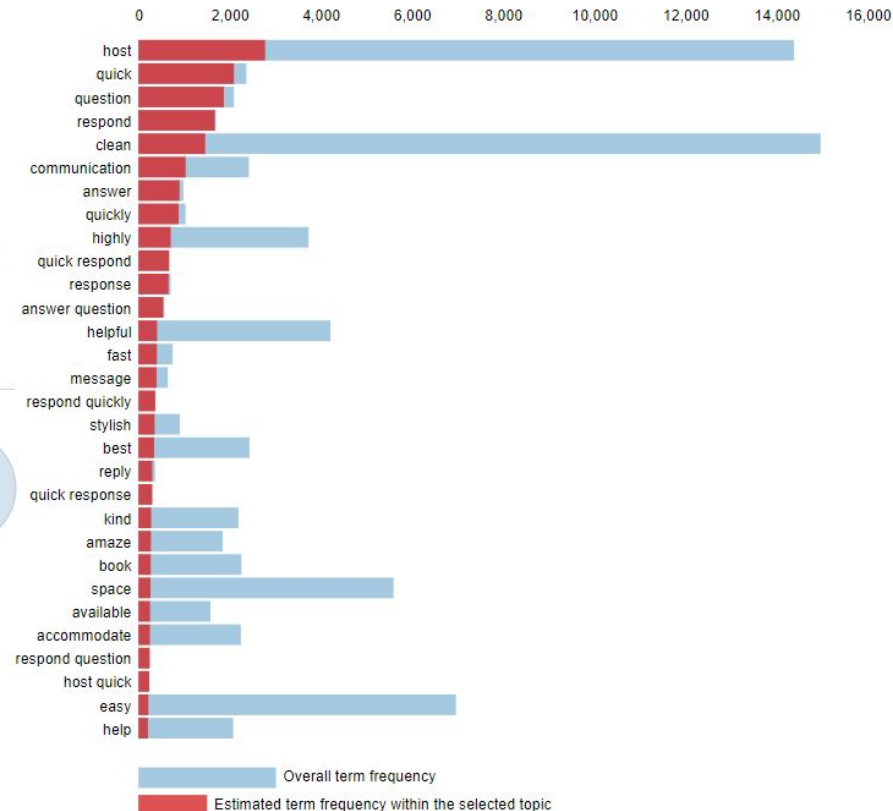
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



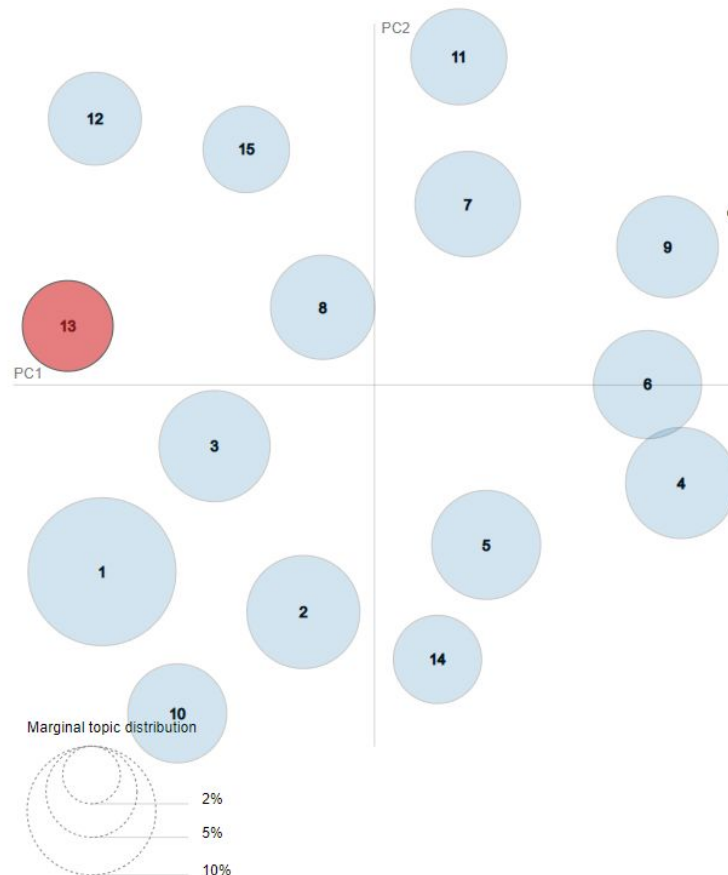
Top-30 Most Relevant Terms for Topic 12 (5.2% of tokens)



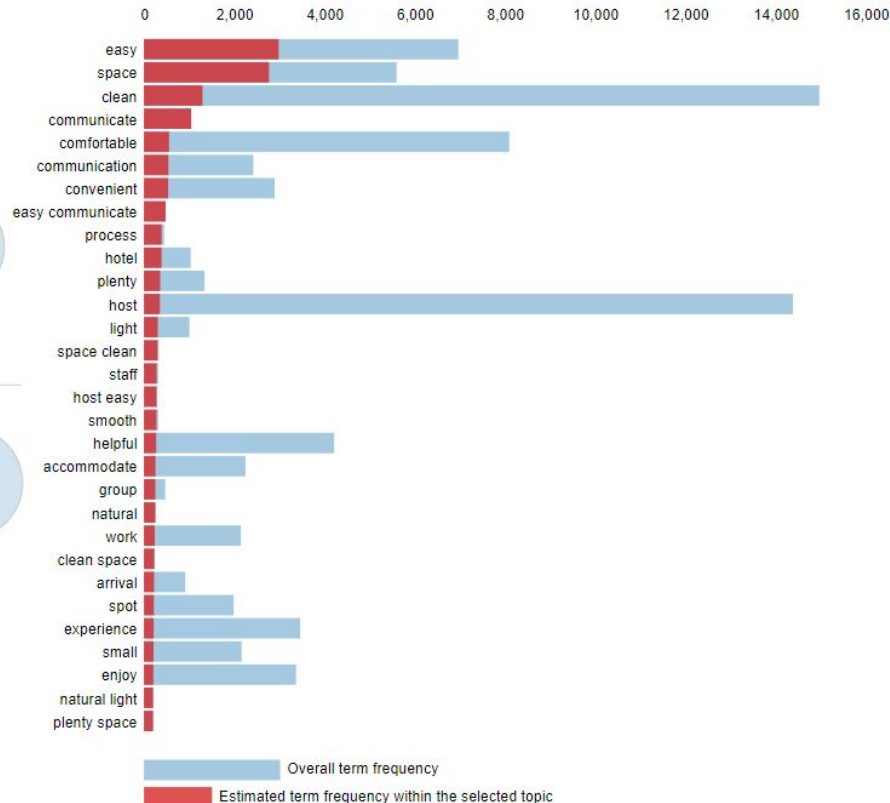
1.  $saliency(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $relevance(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 13 (5% of tokens)



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))]; for topics t; see Chuang et. al (2012)  
 2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

Selected Topic: 14 Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>

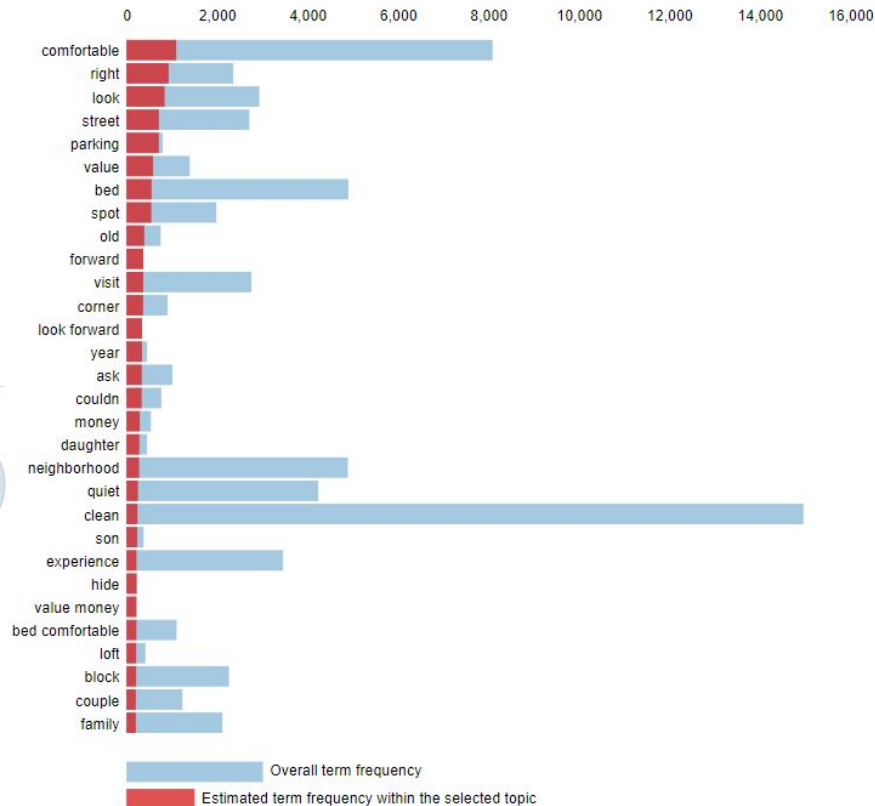
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 14 (4.7% of tokens)



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) =  $\lambda$  \* p(w | t) + (1 -  $\lambda$ ) \* p(w | t)/p(w); see Sievert & Shirley (2014)

Selected Topic:

Slide to adjust relevance metric:(2)

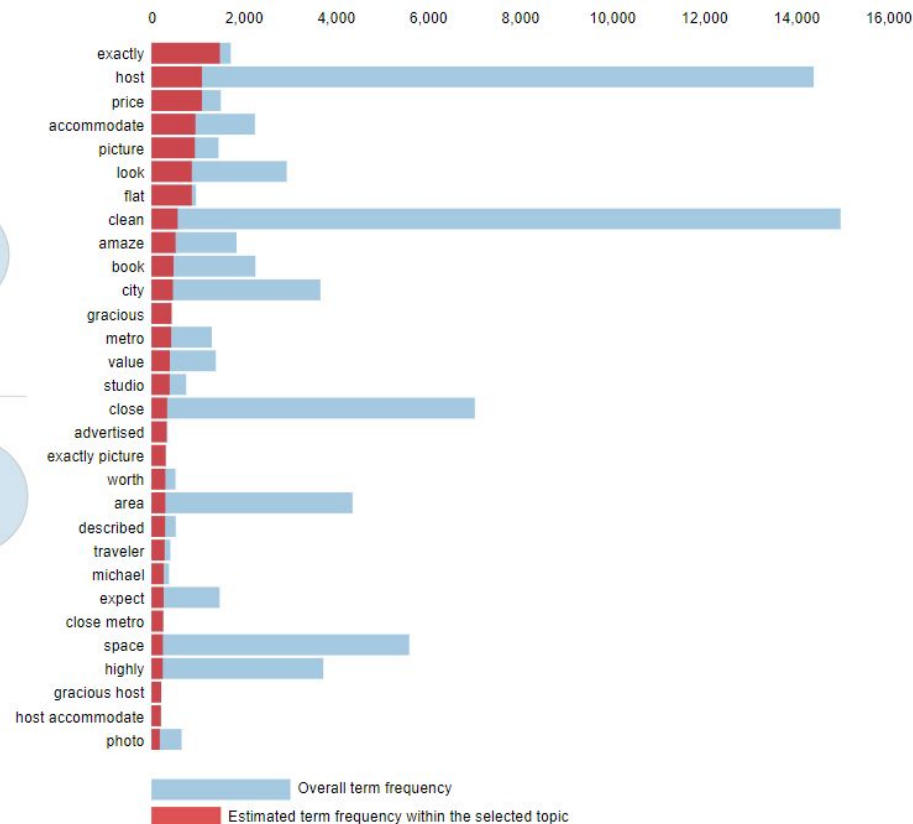
$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 15 (4.5% of tokens)



1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))]] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

## **Negative Topics:**

- Topic 1 - Inaccurate listings
- Topic 2 - Check-in/out
- Topic 3 - Bed/Bathroom
- Topic 4 - Dirtiness and smell
- Topic 5 - Uncomfortable sleep conditions
- Topic 6 - Location
- Topic 7 - Poor house maintenance
- Topic 8 - Noise
- Topic 9 - Hot Water/Heater
- Topic 10 - Dirtiness
- Topic 11 - Location
- Topic 12 - Location general
- Topic 13 - Value

## **Positive Topics Topics:**