



# NLP Project Hotel Review Analysis

METIS  
Project V



*"Great Stay; awesome location"*

★★★★★ Reviewed August 12, 2014

Stayed here for 2 nights for a for a quick trip to SF. Got a great deal off

Thank you for the wonderful review of your recent stay with us! We're thrilled to hear of your friendly interactions with our staff and concierge! We work very hard to ensure that all guests receive top notch customer service while staying with us. We're happy to hear that our guest rooms could meet your needs and that you enjoyed the...



Mustard333  
New York, NY

*"So pretty, but..."*

★★★★☆ Reviewed May 28, 2014

A pretty hotel with enthusiastic staff that just need a little more training. The menu in the dining room was innovative and very good. The upstairs, outdoor bar was a very nice addition. Beautiful, comfortable, well designed hotel/rooms. A very good place to stay in Cincinnati.

What should  
hotels watch out  
for to improve  
their business?

# Goal



Group relevant  
**keywords**



**Suggest** key findings  
for improvement

# Workflow

Preprocess

Data  
NLTK  
Tokenizer



CountVectorize  
ScikitLearn

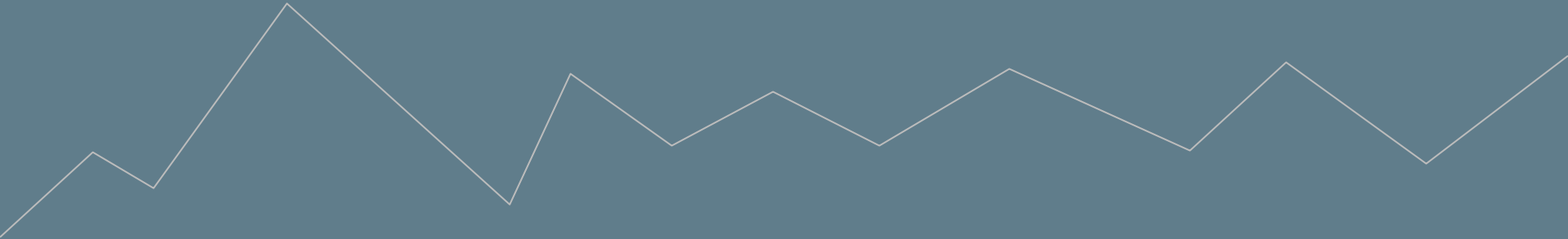


LDA



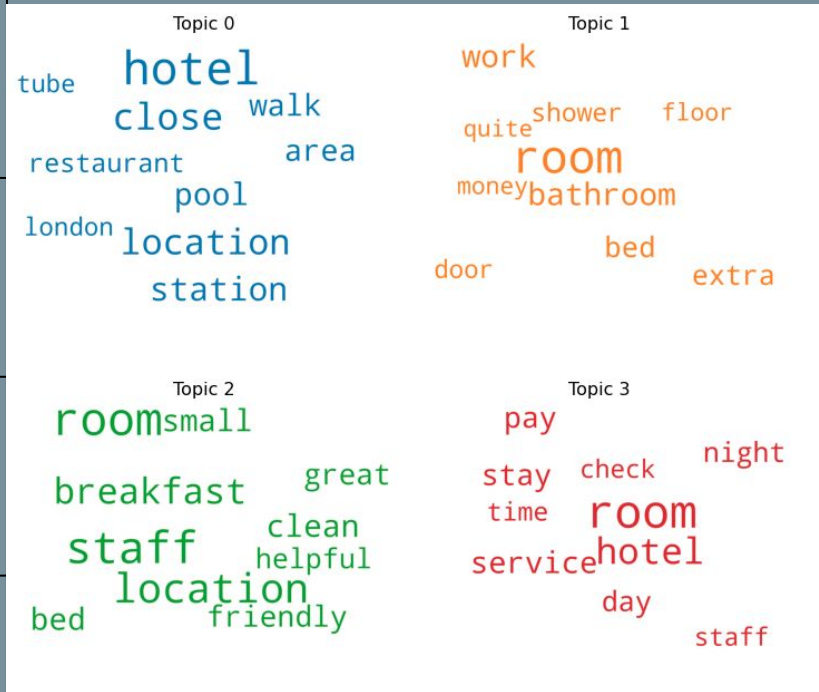
Interpret Findings

WordCloud  
pyLDAvis



# LDA Topic Modeling

	<b>Topic 0:</b> Hotel Surroundings	<b>Topic 1:</b> Hotel Amenities	<b>Topic 2:</b> Hotel Service	<b>Topic 3:</b> Hotel Check-In
<b>1st</b>	Hotel	Room	Room	Room
<b>2nd</b>	Close	Work	Staff	Hotel
<b>3rd</b>	Location	Bathroom	Location	Service



# pyLDavis Visualization

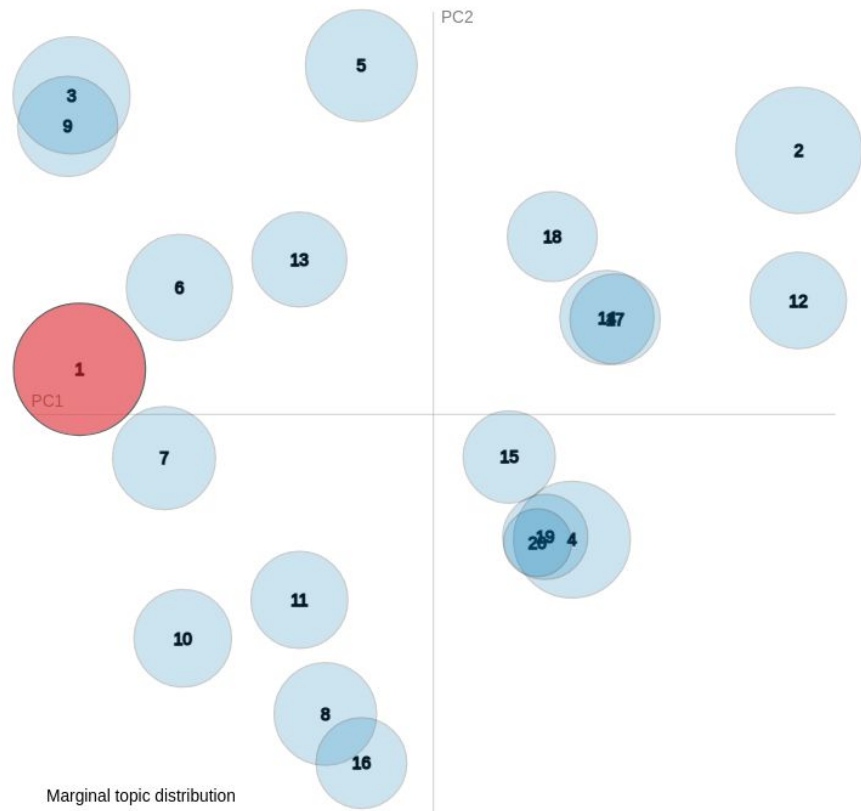
Selected Topic:

Slide to adjust relevance metric:(2)

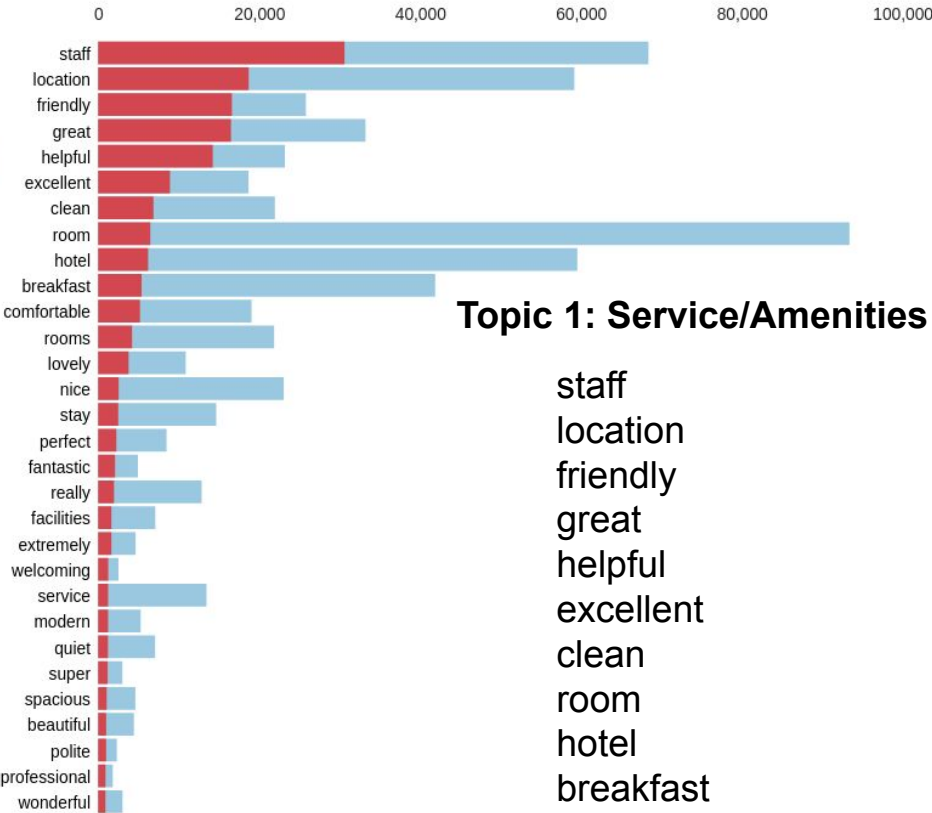
$\lambda = 1$

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 1 (8.4% of tokens)



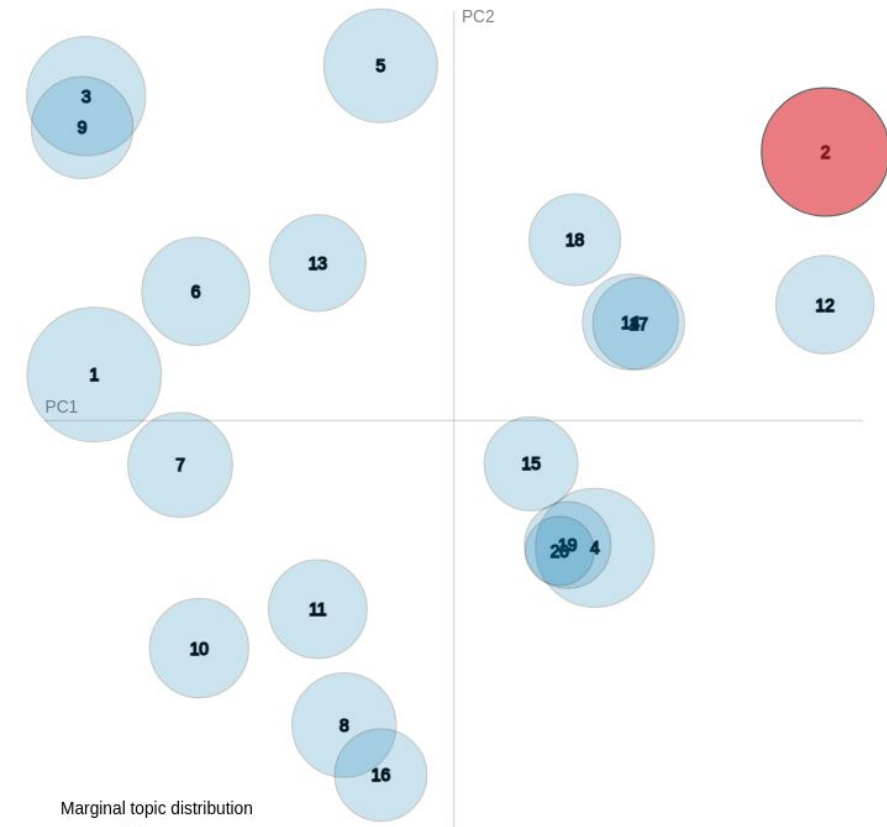
Selected Topic:

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

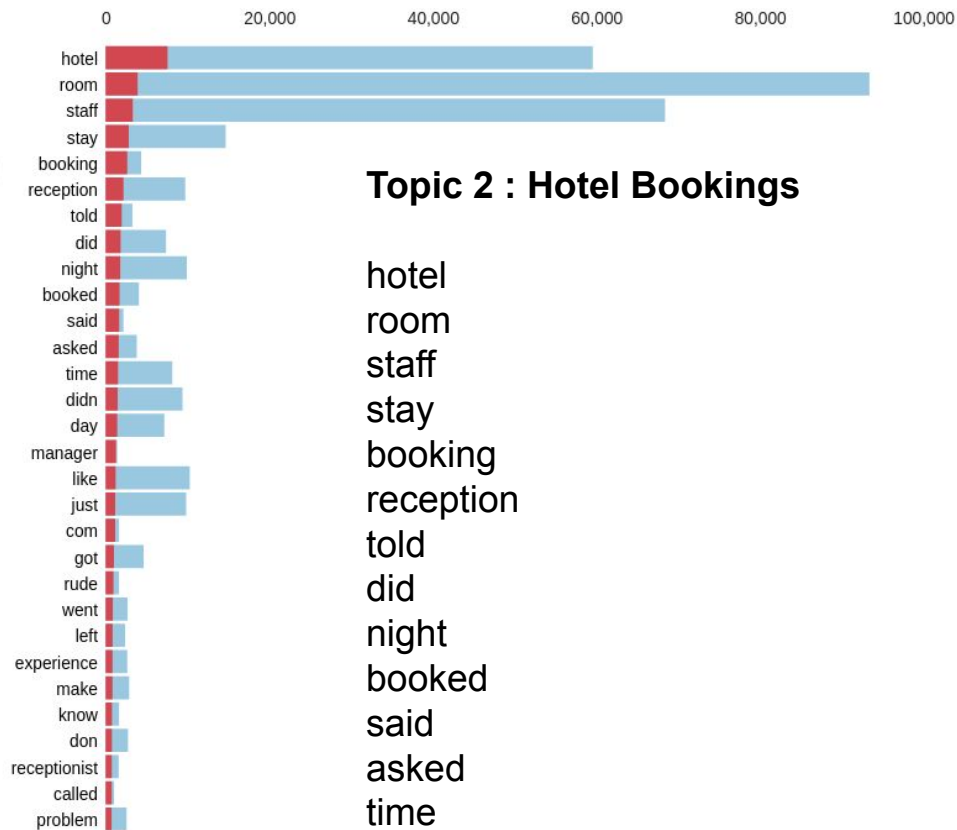
$\lambda = 1$

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)



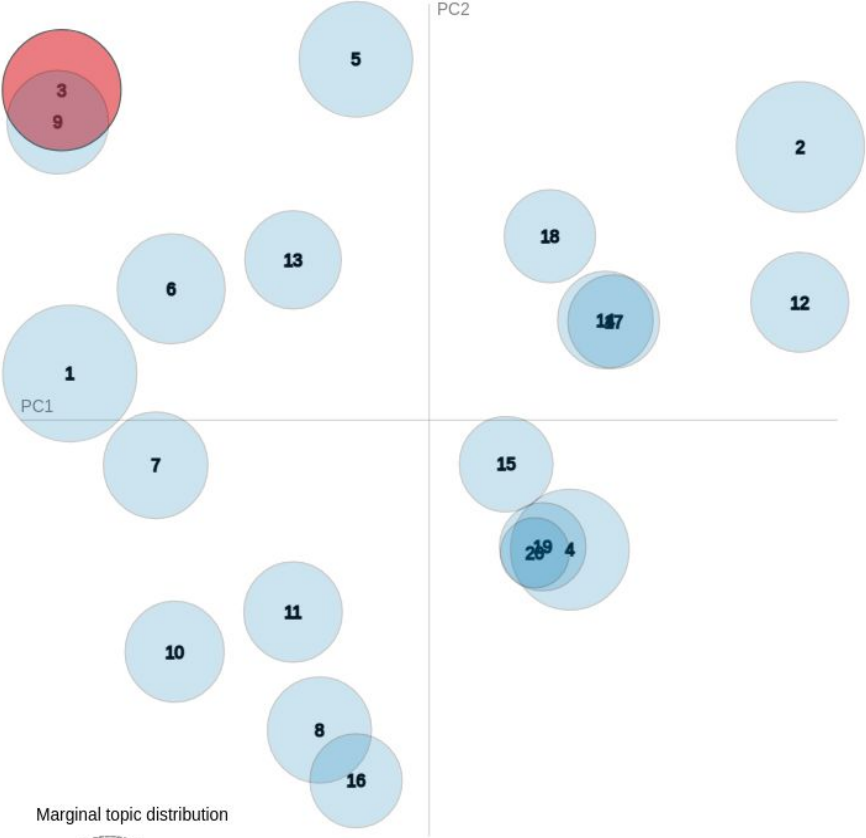
## Topic 2 : Hotel Bookings

hotel  
room  
staff  
stay  
booking  
reception  
told  
did  
night  
booked  
said  
asked  
time

Selected Topic:

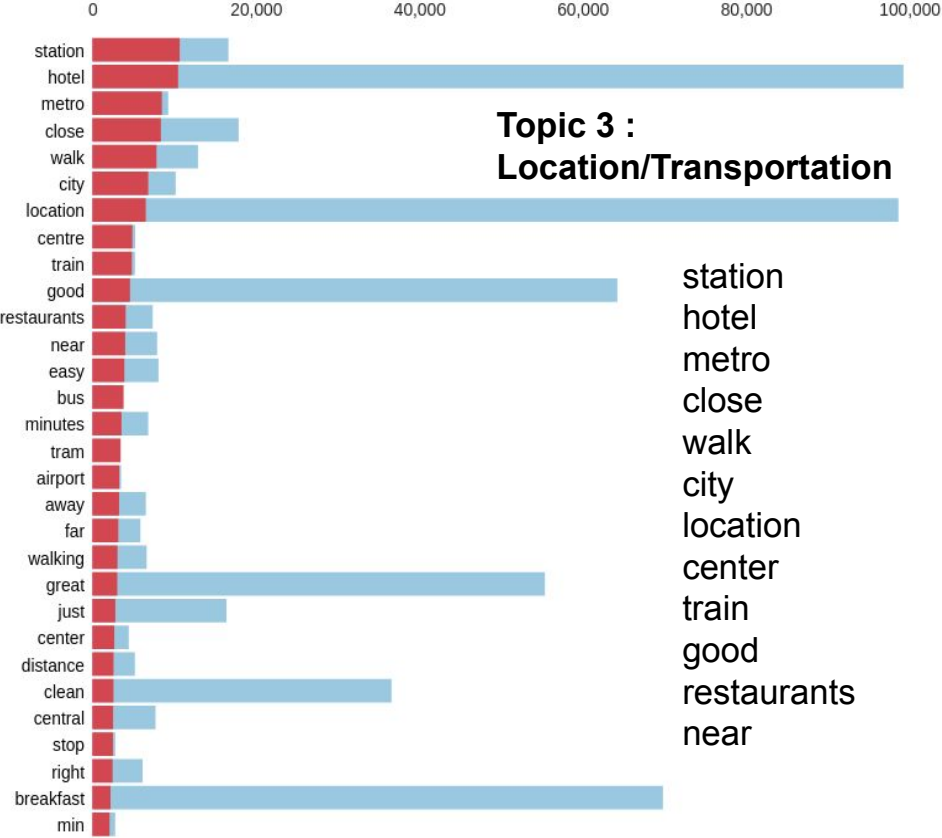
Slide to adjust relevance metric:<sup>(2)</sup>   $\lambda = 1$

Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution

Top-30 Most Relevant Terms for Topic 3 (6.6% of tokens)

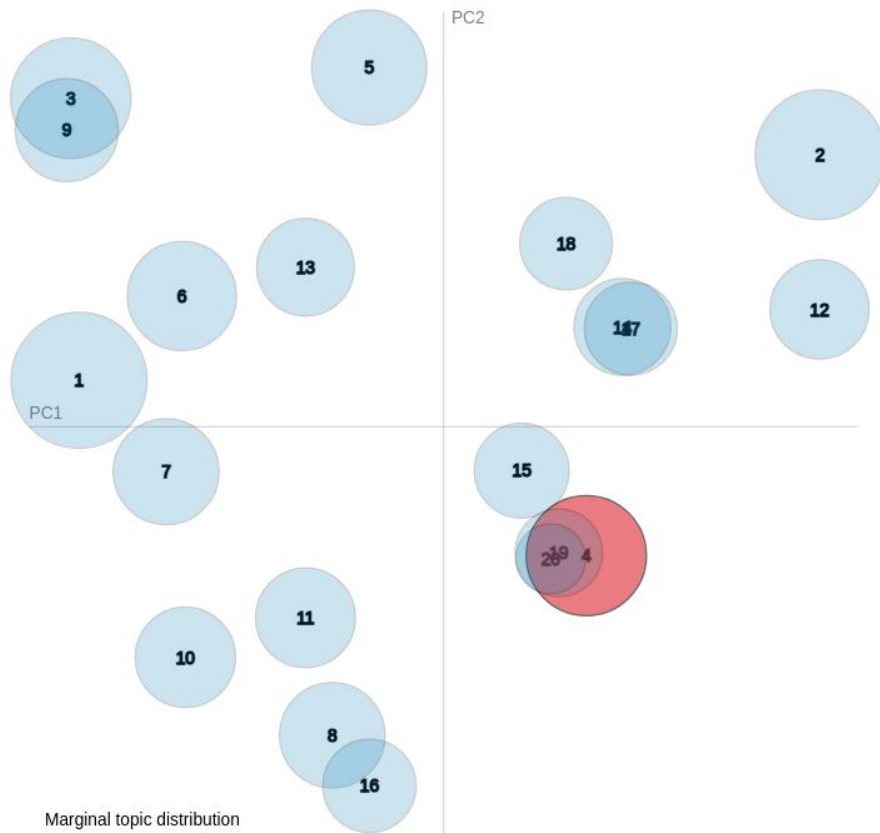




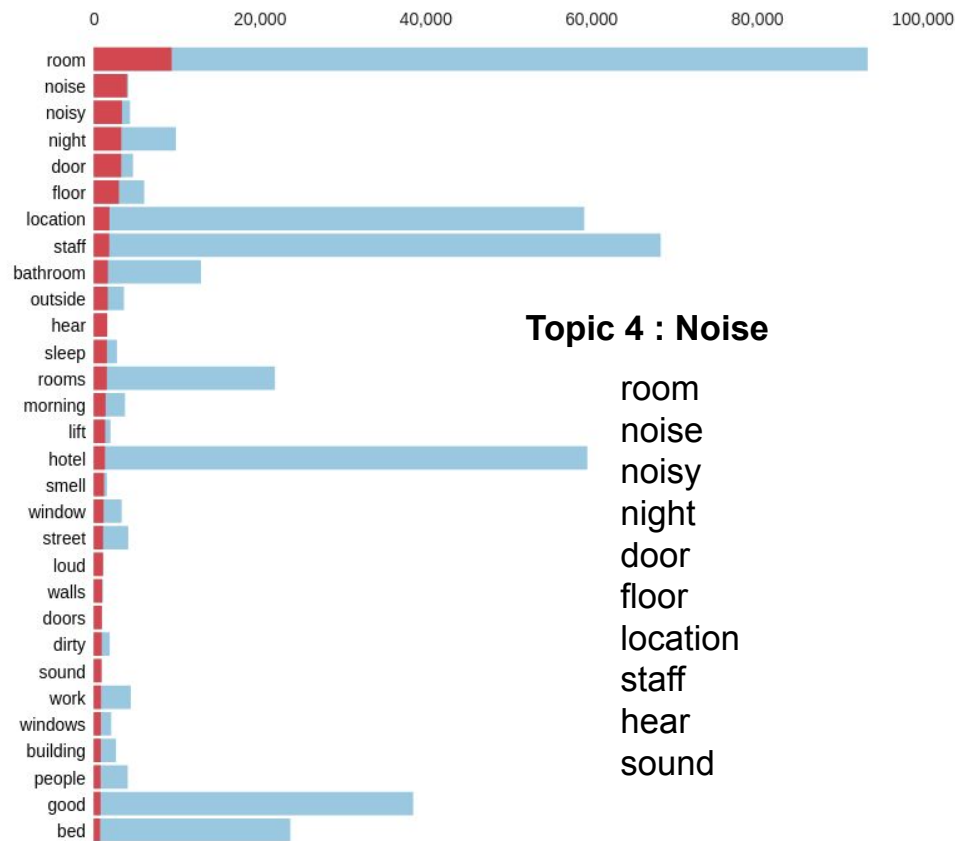
Selected Topic:

Slide to adjust relevance metric: <sup>(2)</sup>   $\lambda = 1$

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 4 (6.6% of tokens)



Selected Topic:

Slide to adjust relevance metric:(2)

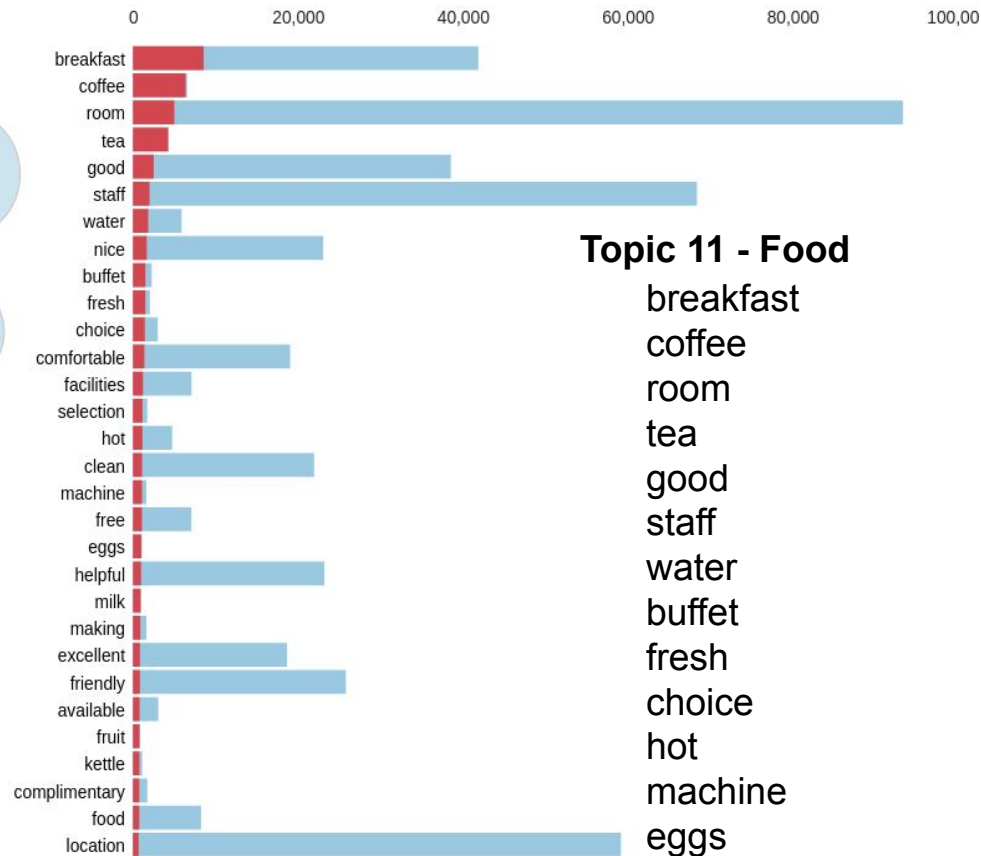
$\lambda = 1$

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 11 (4.6% of tokens)



### Topic 11 - Food

breakfast  
coffee  
room  
tea  
good  
staff  
water  
buffet  
fresh  
choice  
hot  
machine  
eggs

# Overlapping Topics

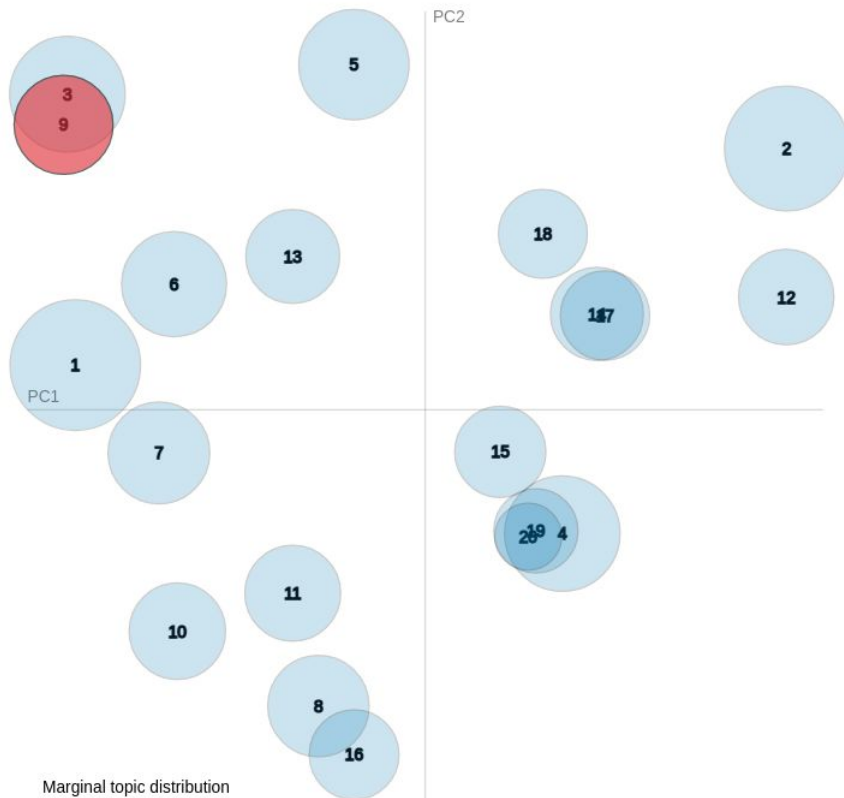
Selected Topic:

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

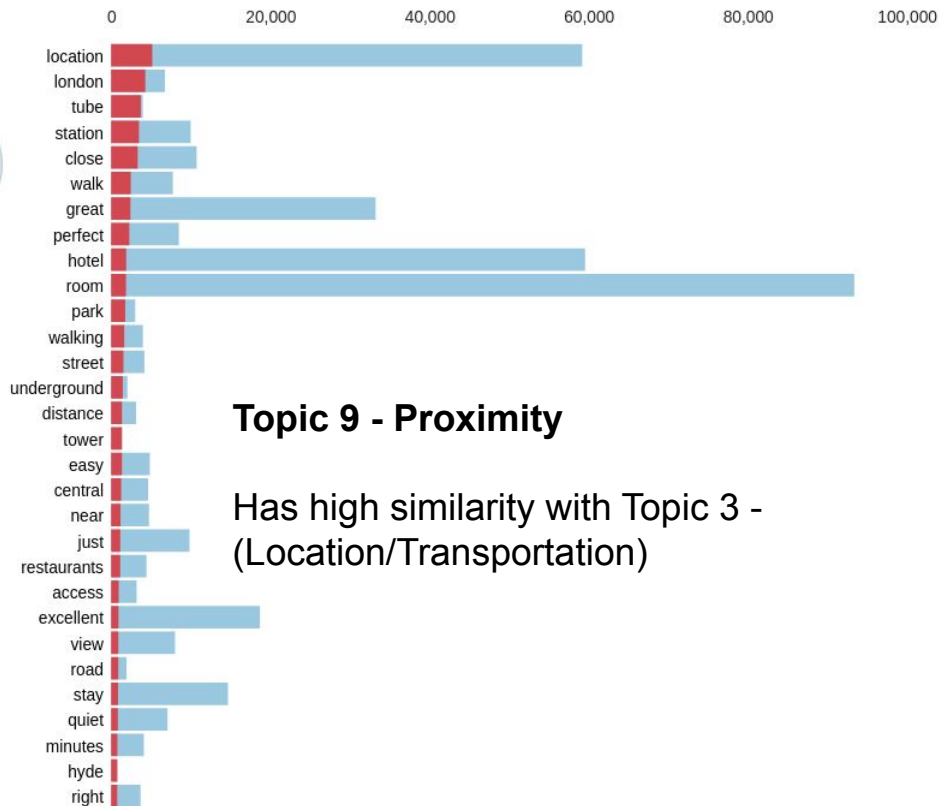
$\lambda = 1$

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 9 (4.9% of tokens)



## Topic 9 - Proximity

Has high similarity with Topic 3 -  
(Location/Transportation)

# Overlapping Topics

Selected Topic:

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

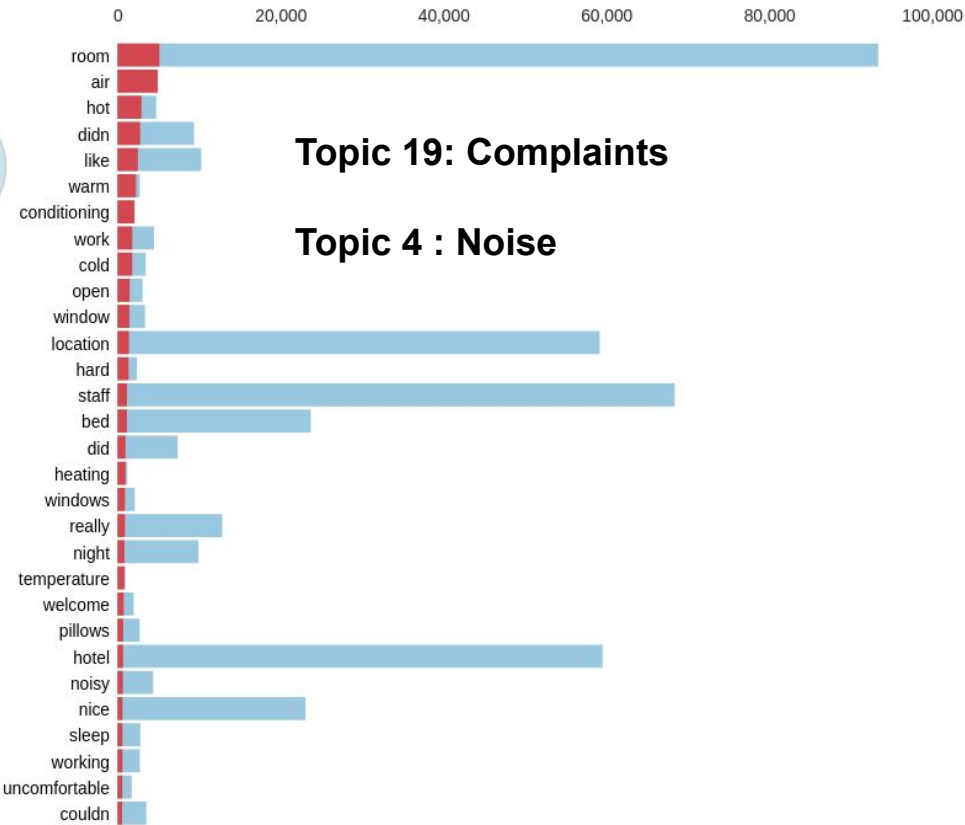
$\lambda = 1$

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 19 (3.5% of tokens)



# Insights/Improvements

- 5 main topics that hotel should focus on
  - **Service, Bookings, Location, Noise, Food**
- Overall a good breakdown of topics
- Overlapping topics are somewhat distinguishable, but share common words

## Improvements

- Coherence score metric to better recommend by topic
  - Create Recommendation System for app users
- 
-

Each bubble = a topic , larger the bubble higher percentage of the number of words in the corpus

Blue bar = overall frequency of each word in corpus

Red bars = estimated # of times a given term was generated by a given topic

# Appendix

## NMF model - Top 15 words

THE TOP 15 WORDS FOR TOPIC 0:

['quite', 'floor', 'big', 'air', 'work', 'window', 'shower', 'noisy', 'little', 'bit', 'view', 'bathroom', 'size', 'small', 'room']

THE TOP 15 WORDS FOR TOPIC 1:

['facility', 'superb', 'quiet', 'wifi', 'ideal', 'expensive', 'comfort', 'fantastic', 'convenient', 'cleanliness', 'price', 'central', 'staff', 'perfect', 'location']

THE TOP 15 WORDS FOR TOPIC 2:

['attentive', 'pleasant', 'super', 'efficient', 'fantastic', 'professional', 'polite', 'really', 'welcome', 'reception', 'lovely', 'extremely', 'helpful', 'friendly', 'staff']

THE TOP 15 WORDS FOR TOPIC 3:

['price', 'close', 'near', 'quality', 'station', 'shower', 'really', 'size', 'restaurant', 'wifi', 'food', 'facility', 'money', 'value', 'good']

# LDA Topic Modeling

```
[(0,
 '0.067*hotel" + 0.041*close" + 0.038*location" + 0.032*station" + '
 '0.029*pool" + 0.027*area" + 0.027*walk" + 0.019*restaurant" + '
 '0.018*london" + 0.017*tube''),
 (1,
 '0.080*room" + 0.028*work" + 0.028*bathroom" + 0.023*bed" + '
 '0.022*extra" + 0.019*shower" + 0.018*door" + 0.016*floor" + '
 '0.016*money" + 0.015*quite''),
 (2,
 '0.088*room" + 0.082*staff" + 0.069*location" + 0.052*breakfast" + '
 '0.038*clean" + 0.037*bed" + 0.037*small" + 0.034*friendly" + '
 '0.034*great" + 0.031*helpful''),
 (3,
 '0.045*room" + 0.026*hotel" + 0.018*service" + 0.016*stay" + 0.016*pay" '
 '+ 0.015*day" + 0.015*night" + 0.013*time" + 0.012*check" + '
 '0.012*staff'')]
```

Topic 0 - Hotel Surroundings

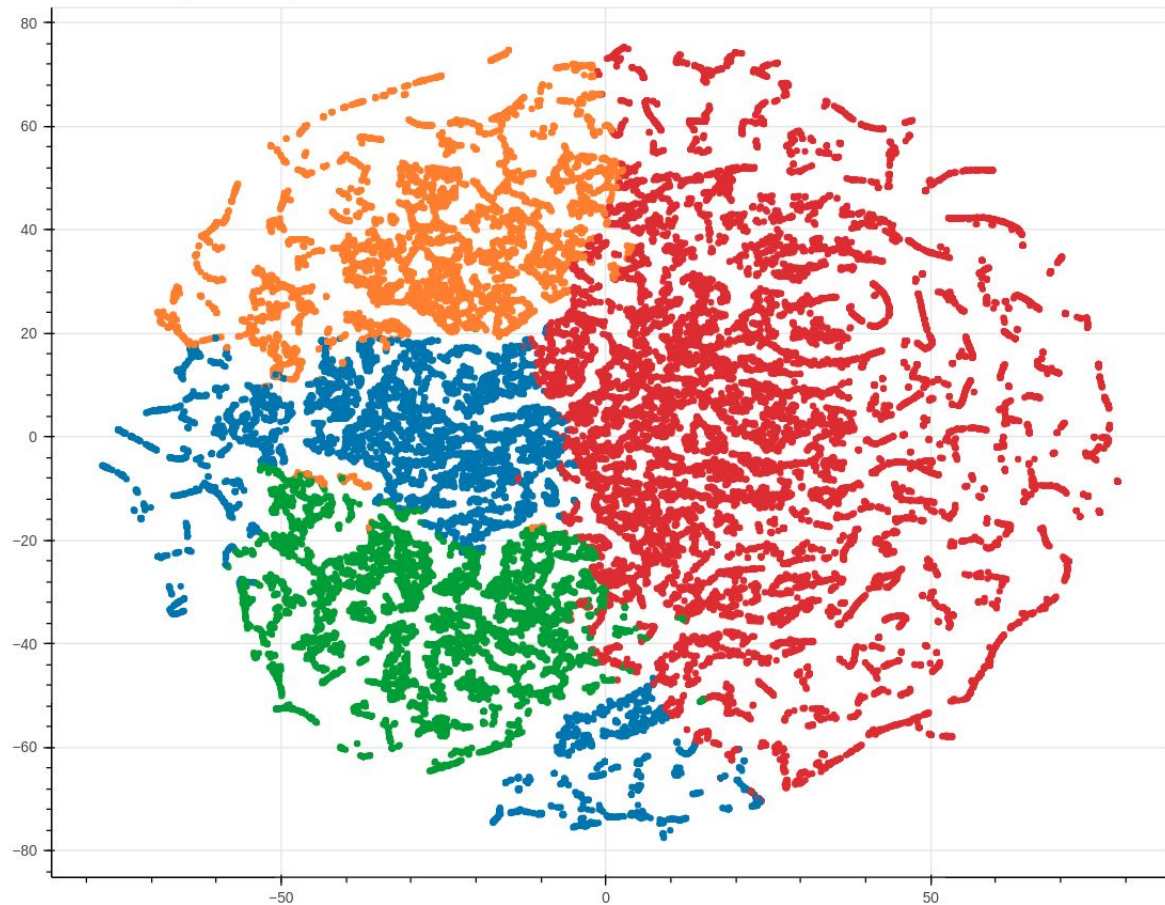
Topic 1 - Hotel Amenities

Topic 2 - Hotel Service

Topic 3 - Hotel Check-In



t-SNE Clustering of 4 LDA Topics



**Red = Topic 3**

**Green = Topic 2**

**Orange = Topic 1**

**Blue = Topic 0**