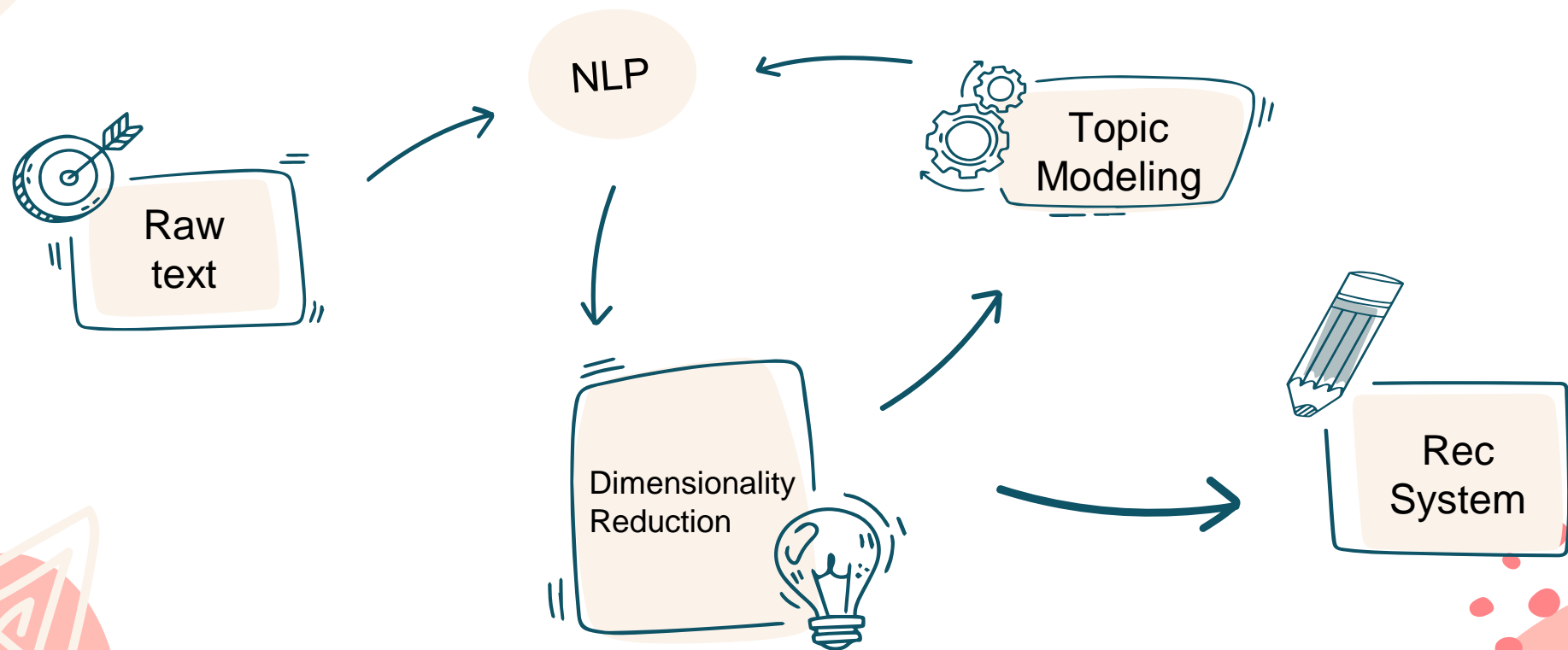# Airbnb Reviews

NLP/Unsupervised Learning

# Objective

Applications:

1.  Topic modeling with reviews can help Airbnb improve customer experience

2.  The reviews recommendation system let's Airbnb find listings that have the most similar reviews

# Workflow



Raw text → NLP → Dimensionality Reduction → Topic Modeling → Rec System

# Methodology

## Data

- Inside Airbnb
- NYC 2009-2021
- 700k data points

## Tools

- NLTK
- Gensim
- SVD
- pyLDAvis
- sklearn

## Models

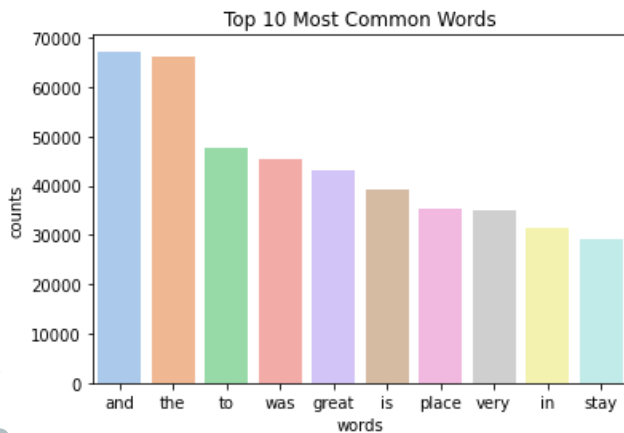- CoreX
- LDA
- NMF

## Recommendation System

- Content-based
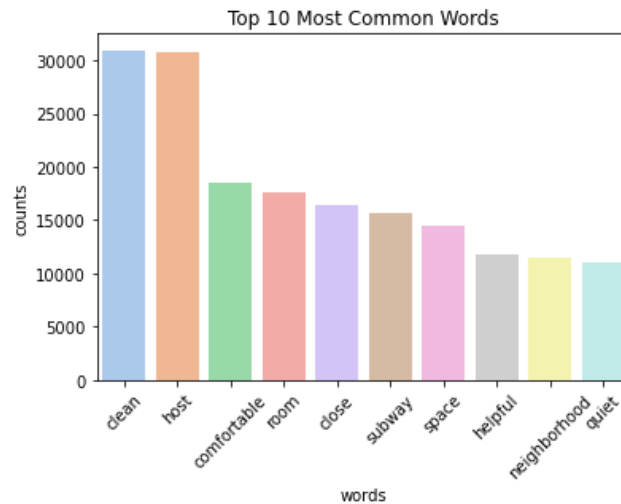
# Text Preprocessing

## Before

48

Average length of document



## After

16

Average length of document

# Topic Modeling

o  NMF/LDA – 3 topics
o  **CoreX** – 5 topics

Anchor words: interior, distance, issues, host, clean

Location — subway, restaurants, distance, walk

Booking Logistics — host, automated, cancel, posting

Host Communications — new, questions, comfortable, available

**Rooms** — **kitchen, bathroom, bed, room**

**Cleanliness** — **clean, sparkling, tidy, spotlessly**

# Recommendation System

○ Tf-IDF → SVD → Rec System

Test review: I really like it she was nice and the places was organize and clean
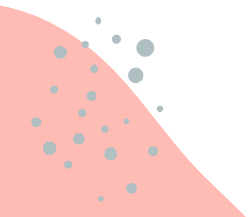
## Airbnb Similar Reviews

Search

Submit

# Future Work

o Obtain more information from Airbnb for further analysis

o Build out recommendation system to include more components

# THANK YOU!

Does anyone have any questions?

# Appendix



LDA 4 topics

LDA 6 topics