

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Understanding and Using Topic Modeling

Using inferred document clusters

Philip Robinson

Presented to OpsLab
NASA - Jet Propulsion Lab

August 24, 2018



Jet Propulsion Laboratory
California Institute of Technology

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Introduction - Philip (1762)

Computer Science MSc at Oregon Health and Science University.



- Whale vocalization
- Light pollution
- Modern cryptography
- ★ Information retrieval



Jet Propulsion Laboratory
California Institute of Technology

Presentation Overview

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

1 Introduction

2 What is Topic Modeling

3 Why Topic Model

4 A Generative Model

5 Understanding Model Space

6 Details

7 Notes

8 Conclusion



Jet Propulsion Laboratory
California Institute of Technology

What is Topic Modeling

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Topic modeling is a text processing technique for automatically grouping documents by topics. This is usually used as a strategy to describe documents in low dimensional space or an exploratory tool for document collections.



Jet Propulsion Laboratory
California Institute of Technology

Examples

In practice, this requires many more documents

The Tourist huddles in the station While slowly night gives way to dawn; He finds a certain fascination In knowing all the trains are gone.

- Food
- Travel
- Time

The Governess up in the attic Attempts to make a cup of tea; Her mind grows daily more erratic From cold and hunger and ennui.

From this annotation we know that Document 2 and 3 are about Food and Time

The Journalist surveys the slaughter, The best in years without a doubt; He pours himself a gin and water and wonders how it came about.



Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

What can I solve?

- Find similar document pairs
- Cluster documents into groups with similar content
- Find relevant documents to a query or user's interests
- Explore shape of document collection
- Can be extended to address many other problems



Jet Propulsion Laboratory
California Institute of Technology

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Applications of Topic Modeling

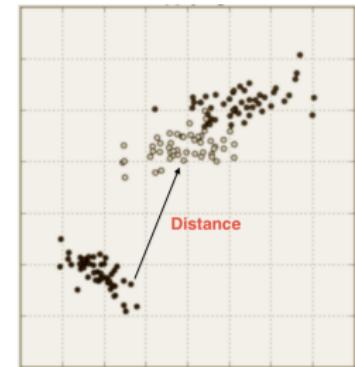
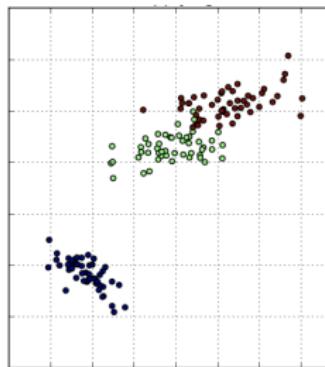
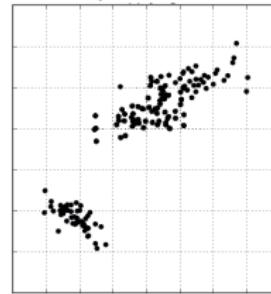
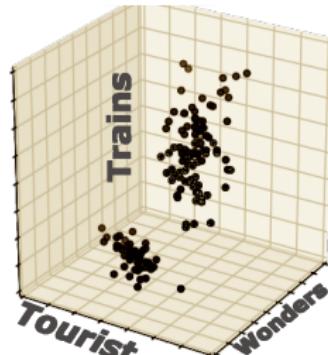


Figure: Exploratory
Data Analysis

Figure: Cluster Points

Figure:
Dimmentionality
Reduction



Jet Propulsion Laboratory
California Institute of Technology

Latent Dirichlet Allocation (LDA)

A generative model

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

- Assume/Generalize how documents could have been generated
- Fit parameters that describe generalization
- Ask questions about the generalization in relation to documents
- Ask questions about documents in relation to the generalization

Topic modeling generalizes how a document is generated by claiming that words come from topics, and documents have multiple topics. Note that this ignores sentence structure, entities, authorship, or other things we may care about.



Jet Propulsion Laboratory
California Institute of Technology

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

- Represent document as Bag-of-Words¹
- Model/Fit topics as mixture of words
- Documents are projected into topic space
- Study relationships between document projections

Topic Modeling

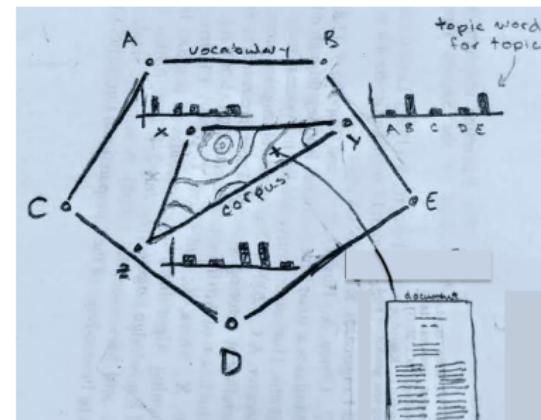


Figure: Latent Dirichlet Allocation

$$T(x) = \text{Project } x \text{ into topic-space}$$

$$\text{Measure} = \text{Distance}(T(\text{Doc 1}), T(\text{Doc 2}))$$

¹equivalent to multinomial over the vocabulary



Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Dirichlet Distribution

In this case the Topic-Space is our Dirichlet Distribution

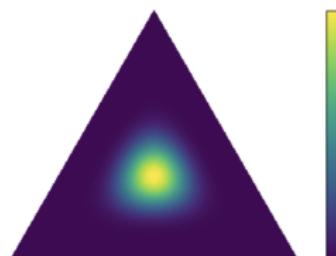


Figure:
Non-Informative

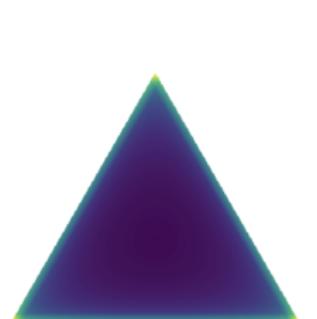


Figure: Little In
Common

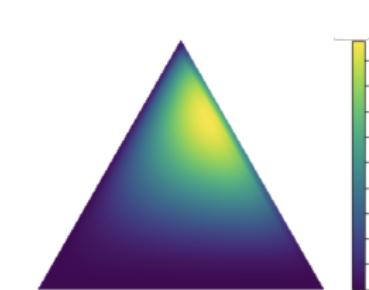


Figure: Shared Topics



Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Using this in Python²

```
from nltk.corpus import brown

data = []

for fileid in brown.fileids():
    document = ' '.join(brown.words(fileid))
    data.append(document)

NO_DOCUMENTS = len(data)
print(NO_DOCUMENTS)
print(data[:5])

from sklearn.decomposition import LatentDirichletAllocation
from sklearn.feature_extraction.text import CountVectorizer

NUM_TOPICS = 10

vectorizer = CountVectorizer(min_df=5, max_df=0.9,
                             stop_words='english', lowercase=True,
                             token_pattern='[a-zA-Z\\-][a-zA-Z\\-]{2,}')
data_vectorized = vectorizer.fit_transform(data)

# Build a Latent Dirichlet Allocation Model
lda_model = LatentDirichletAllocation(n_topics=NUM_TOPICS, max_iter=10, learning_method='online')
lda_Z = lda_model.fit_transform(data_vectorized)

text = "The economy is working better than ever"
x = lda_model.transform(vectorizer.transform([text]))[0]
```

²example in scikit-learn, I used gensim



Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Concerns about Topics

- Topics do not have labels
- Topics are not easily human interpretable
- In traditional LDA n-grams are not represented



Visualizing High Dimensional Space³

Introduction

What is Topic
Modeling

Why Topic Model

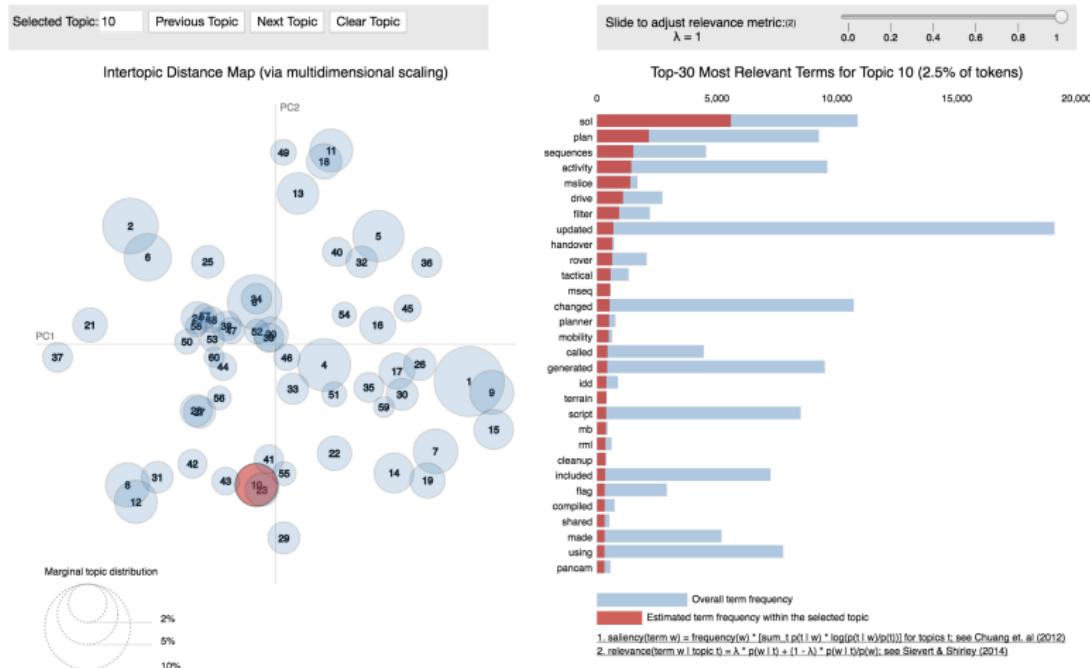
A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion



³used pyLDAvis



Jet Propulsion Laboratory
California Institute of Technology

Looking at Topics

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Breakout to Jupyter



Jet Propulsion Laboratory
California Institute of Technology

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Evaluation Notes

As this is a bayesian machine learning approach we will need to evaluate model fit.

- perplexity
- coherence
- visualization
- predictive power



Jet Propulsion Laboratory
California Institute of Technology

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Preprocessing Notes

As this model generates conditional probabilities on observed words, we need a clean/normalize vocabulary

- Lowercase corpus
- Tokenization & (Stemming — Lemmatization)
- Removing symbols
- Remove highly common words
- Remove extremely rare words



Jet Propulsion Laboratory
California Institute of Technology

Introduction

What is Topic
Modeling

Why Topic Model

A Generative
Model

Understanding
Model Space

Details

Notes

Conclusion

Conclusion & Questions

I used this technique during my internship to automatically assign tickets to subject matter experts, for the Office of Safety and Mission Success (5x)⁴

⁴1762 acts as data science support to JPL

