# LDA with Gibbs Sampling

Group 3

# Introduction

# Evolution of topic models:
# TF-IDF → LSA → *pLSA* → LDA

TF-IDF is word-based, not topic-based

LSA leverages TF-IDF to identify weights of words, which are then compressed into set number of topics

Latent Dirichlet Allocation - probabilistic approach

# LDA is a probabilistic topic model, departing from LSA and TF-IDF

|  | **LSA** | **LDA** (Blei et al. (2003)) |
|---|---|---|
| *Approach* | Matrix factorization | Probabilistic |
| *Generative?* | No, not able to apply to "unseen" documents | Yes, can be applied to "unseen" documents |
| *Algorithm* | Harder to implement | Easier to implement |

LSA & LDA shared assumptions:
1) bag of words
2) order of documents is not significant

# LDA origins

The paper "Latent Dirichlet Allocation" was published by Blei, Ng and Jordan in the Journal of Machine Learning Research in 2003

LDA was independently invented for use in population genetics research by Pritchard, Stephens and Donnelly in 2000

# What is LDA?

- **Latent Dirichlet Allocation**
  - Words are the only observable variables, all others are **latent** variables
  - Leverages **Dirichlet** distributions
  - **Allocates** the words of the document to different topics

- Generative Statistical Model for Topic Modeling
  - Imagine how the documents were created and reverse engineer generation

- Documents contain multiple topics, but probably not all of them

# Generative Process

Each **topic** is represented as a distribution over a fixed vocabulary

 e.g. a *genetics* topic would have a high probability of containing words about genetics

Each **document** has a topic distribution

 e.g. an article about genetic data analysis would have a high probability for the topics *genetics* and *data analysis*

Each **word** in a document is chosen from a topic

 - the topic *genetics* might be chosen with a high probability from the topic distribution for the document

 - the word *gene* might be chosen from the topic *genetics* with high probability

# Generative Process

# What kind of distribution?

**Dirichlet Distribution**

Distribution over multinomial distributions

P(word = mouse) = (½, ½, 0)

# Concentration parameter

# Inference

To "generate" a new document, we need to know:

1. **Topic distribution** for the document
2. **Word distributions** for each topic

We don't know the distributions, so we have to **infer** them from training data

That's computationally expensive!

Sample from the distribution to iteratively approximate the values

# How do we sample?

**Gibbs Sampling**

Approximate joint distributions for latent variables to sample from

Remove one value for a latent variable, then calculate new joint distribution conditioned on other values, then randomly sample from that distribution

E.g. For each word, unassign a topic, compute a new joint distribution for that word and each topic based on all other words with topic assignments in the document, choose a new topic assignment from the distribution

Update priors based on observations

# High Level Overview of Algorithm

Input:

    Set of documents, made up of words

    Number of topics to find

Learning:

    Initialize topic distributions and topic-word distributions (usually randomly)

    Using a Gibbs Sampler, iteratively sample and update

    Update priors until stop condition

Output:

    Topic distributions and topic-word distributions

# Algorithm

# LDA Generation Plate Notation

# Generative Process: the assumption

If we have a document of a certain length:

And we want it to be 60% about fashion

and 40% about business:

**Document**
word word word word word word word
word word word word word word word
word word word word word word word
word word word word word word word

# Generative Process: the assumption

If we have a list of fashion words and a list of business words:

We could generate a document using 60% fashion words and 40% business words:

| **Fashion** | **Business** |
| --- | --- |
| clothes shoes | market value |
| design style | company profit |
| vogue trend | trend merger |
| popular hip | revenue sell |

**Document**

trend hip clothes clothes profit hip merger hip market clothes trend popular trend style trend hip design profit company style

$$p(\theta, \mathbf{z} \,|\, \mathbf{w}, \alpha, \beta) = \frac{p(\theta, \mathbf{z}, \mathbf{w} \,|\, \alpha, \beta)}{p(\mathbf{w} \,|\, \alpha, \beta)}.$$

# Gibbs Sampling

- Markov Chain Monte Carlo (MCMC) algorithm
- Sample from distributions with two or more dimensions
- When the conditional probabilities can be calculated, estimate joint probabilities.

# Simple 2-D Gibbs example

|           | mouse | horse | money | market |
|-----------|:-----:|:-----:|:-----:|:------:|
| finance   |   0   |   1   |   4   |   5    |
| computers |   5   |   0   |   2   |   1    |
| animals   |   3   |   3   |   0   |   2    |

|           | mouse | horse | money | market |
|-----------|:-----:|:-----:|:-----:|:------:|
| finance   | 0 | $\frac{1}{10}$ | $\frac{2}{5}$ | $\frac{1}{2}$ |
| computers | $\frac{5}{8}$ | 0 | $\frac{1}{4}$ | $\frac{1}{8}$ |
| animals   | $\frac{3}{8}$ | $\frac{3}{8}$ | 0 | $\frac{1}{4}$ |

$$P(word|topic)$$

|           | mouse | horse | money | market |
|-----------|:-----:|:-----:|:-----:|:------:|
| finance   | 0 | $\frac{1}{4}$ | $\frac{2}{3}$ | $\frac{5}{8}$ |
| computers | $\frac{5}{8}$ | 0 | $\frac{1}{3}$ | $\frac{1}{8}$ |
| animals   | $\frac{3}{8}$ | $\frac{3}{4}$ | 0 | $\frac{1}{4}$ |

$$P(topic|word)$$

|  | mouse | horse | money | market |
|---|---|---|---|---|
| finance | $0$ | $\frac{1}{4}$ | $\frac{2}{3}$ | $\frac{5}{8}$ |
| computers | $\frac{5}{8}$ | $0$ | $\frac{1}{3}$ | $\frac{1}{8}$ |
| animals | $\frac{3}{8}$ | $\frac{3}{4}$ | $0$ | $\frac{1}{4}$ |

$$P(topic|word)$$

|  | mouse | horse | money | market |
|---|---|---|---|---|
| finance | 0 | $\frac{1}{10}$ | $\frac{2}{5}$ | $\frac{1}{2}$ |
| computers | $\frac{5}{8}$ | 0 | $\frac{1}{4}$ | $\frac{1}{8}$ |
| animals | $\frac{3}{8}$ | $\frac{3}{8}$ | 0 | $\frac{1}{4}$ |

$$P(word|topic)$$

|           | mouse         | horse         | money         | market        |
| --------- | ------------- | ------------- | ------------- | ------------- |
| finance   | $0$           | $\frac{1}{4}$ | $\frac{2}{3}$ | $\frac{5}{8}$ |
| computers | $\frac{5}{8}$ | $0$           | $\frac{1}{3}$ | $\frac{1}{8}$ |
| animals   | $\frac{3}{8}$ | $\frac{3}{4}$ | $0$           | $\frac{1}{4}$ |

$$P(topic|word)$$

|          | mouse         | horse          | money         | market        |
|----------|---------------|----------------|---------------|---------------|
| finance  | 0             | $\frac{1}{10}$ | $\frac{2}{5}$ | $\frac{1}{2}$ |
| computers| $\frac{5}{8}$ | 0              | $\frac{1}{4}$ | $\frac{1}{8}$ |
| animals  | $\frac{3}{8}$ | $\frac{3}{8}$  | 0             | $\frac{1}{4}$ |

$$P(word|topic)$$

1. Randomly initialize each $x_i$

2. For $t = 1, ..., T$:

    2.1   $x_1^{t+1} \sim p(x_1 | x_2^{(t)}, x_3^{(t)}, ..., x_m^{(t)})$

    2.2   $x_2^{t+1} \sim p(x_2 | x_1^{(t+1)}, x_3^{(t)}, ..., x_m^{(t)})$

    2.m   $x_m^{t+1} \sim p(x_m | x_1^{(t+1)}, x_2^{(t+1)}, ..., x_{m-1}^{(t+1)})$

# Inference

For all the words in every document, start out by assigning each one to a topic randomly:

| Document |
|---|
| **Document** <br> cat dog animal dog animal cat cat |
| **Document** <br> apple pie ingredient apple pie flour dough |
| **Document** <br> market analyst invest invest price market |

| | |
|---|---|
| Topic 1 | cat cat  apple apple dough  invest market |
| Topic 2 | dog animal cat  pie pie  market invest |
| Topic 3 | animal dog  ingredient flour  analyst price |

# Inference

Count the number of times each word occurs with each topic...

|  | cat | dog | animal | apple | pie | market | ... |
|---|---|---|---|---|---|---|---|
| Topic 1 | 2 | 0 | 0 | 2 | 0 | 1 | |
| Topic 2 | 1 | 1 | 1 | 0 | 2 | 1 | |
| Topic 3 | 0 | 1 | 0 | 0 | 0 | 0 | |

# Inference

And count the number of times words from each document occur with each topic.

|  | Document 1 words | Document 2 words | Document 3 words |
|---|---|---|---|
| Topic 1 | 2 | 3 | 2 |
| Topic 2 | 3 | 2 | 2 |
| Topic 3 | 2 | 2 | 2 |

# The Algorithm

For every document *d:*

For every word *w* in the document:

For every topic *t*:

$$P(t) \sim \frac{\text{\# of times } w \text{ occurs in that topic } (+\beta)}{\text{\# of times } w \text{ occurs in that topic } + \text{\# of unique words} * \beta} * \frac{\text{\# of words in document in that topic } (+\alpha)}{\text{\# of words } + \text{\# of topics} * \alpha}$$

# The Algorithm

**Document**

cat dog animal dog animal

cat cat

For every document *d*:

For every word *w* in the document:    cat

For every topic *t*:

| Topic 1 | ~~cat~~ ~~cat~~ apple apple dough invest market |
|---|---|

$$P(t) \sim \frac{\text{\# of times } w \text{ occurs in that topic } (+\beta)}{\text{\# of times } w \text{ occurs in that topic } + \text{\# of unique words} * \beta} * \frac{\text{\# of words in document in that topic } (+\alpha)}{\text{\# of words } + \text{\# of topics} * \alpha}$$

$$\frac{1 + \beta}{1 + 12 * \beta} * \frac{1 + \alpha}{12 + 3 * \alpha}$$

=0.05555

# The Algorithm

| Document |
|:---:|
| **Document** |
| cat dog animal dog animal cat cat |

cat

Alpha = 0.5
Beta = 0.1

Topic 1 ~ 0.05555
Topic 2 ~ 0.12963
Topic 3 ~ 0.01543

$P(t = 1 \mid w, d, z) = (0.05555) / (0.05555 + 0.12963 + 0.01543)$
$= 0.27691$
$P(t = 2 \mid \ldots) = 0.64618$
$P(t = 3 \mid \ldots) = 0.07692$

# The Algorithm

Randomly sample from this distribution:

# The Algorithm

Reassign the word to the new topic. Do this over and over and over again.

| Topic 1 | ~~cat~~ cat  apple apple dough  invest market |
| --- | --- |
| Topic 2 | **cat** dog animal cat  pie pie  market invest |
| Topic 3 | animal dog  ingredient flour  analyst price |

# Resulting Text Representation

Documents, which are a bag of words, are represented as a mixture of topics from which words can be sampled from multinomial distributions.

# Performance

# Intrinsic evaluation metrics

Hold-out perplexity

$$\hat{H} = -\frac{1}{m} \log_2 P(w_1, w_2, \ldots, w_m)$$

$$PP = 2^{\hat{H}}$$

Coherence (PMI, NPMI)

$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)}$$

$$\text{npmi}(x; y) = \frac{\text{pmi}(x; y)}{h(x, y)}$$

# Extrinsic evaluation metrics

If your gold standard data are labeled, accuracy can be measured directly by use of downstream algorithms (SVM, clustering, etc).

Multiclass classification: cross-entropy

Clustering: B-cubed, F-measure, etc.

# Computational complexity

Inference is $\theta$(k * |d| * |V|) where k is number of topics, d is the set of documents, V is the vocab.

For a large number of topics, some research suggests that the problem is NP-hard.

# Demo

Go to Jupyter notebook

# Application

# LDA for recommending NYT articles

**User**: distribution of topics they're interested in

**Article**: distribution of topics/words

Adjust topic distributions based on reader preferences
- Add offsets to model topic error, incorporate reading patterns
- Iteratively adjust offsets and then recalculate reader scores

# Grade of Membership (GoM) Models: Population genetics equivalent to LDA



*Image from Raj, Stephens and Pritchard, "fastSTRUCTURE: Variational Inference of Population Structure in Large SNP Data Sets" Genetics (2014)*

populations
=
topics

DNA segments
=
documents

genes
=
words

# Questions?

# References

- https://en.wikipedia.org/wiki/Latent_Dirichlet_allocation
- https://en.wikipedia.org/wiki/Gibbs_sampling
- http://www.cs.columbia.edu/~blei/papers/WangBlei2011.pdf
- https://open.blogs.nytimes.com/2015/08/11/building-the-next-new-york-times-recommendation-engine/
- https://wiseodd.github.io/techblog/2017/09/07/lda-gibbs/
- http://www.cs.columbia.edu/~blei/papers/Blei2012.pdf
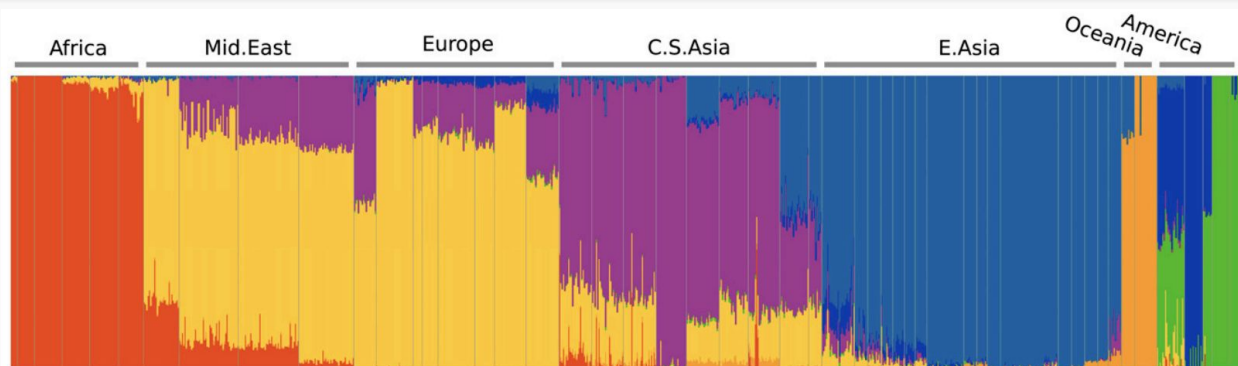- http://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf
- http://www.genetics.org/content/197/2/573
- https://stats.stackexchange.com/questions/244917/what-exactly-is-the-alpha-in-the-dirichlet-distribution
- https://www.youtube.com/watch?v=yK7nN3FcgUs&feature=youtu.be
- https://www.youtube.com/watch?v=FkckgwMHP2s
- https://people.cs.umass.edu/~wallach/publications/wallach09evaluation.pdf
- http://cseweb.ucsd.edu/~dhu/docs/exam09.pdf

# References

- http://www1.icsi.berkeley.edu/Speech/docs/HTKBook3.2/node188_mn.html
- https://en.wikipedia.org/wiki/Pointwise_mutual_information
- https://towardsdatascience.com/dirichlet-distribution-a82ab942a879

# Credits

Aidan: Generative process, inference, algorithm slides

Elijah: Initial intro, performance, demo, presentation

Julia: Intro, Applications - general NLP & applications beyond linguistics

Kevin: Gibbs algorithm example and algorithm overview slides

Zoe: Intro/General Overview, NLP applications, further applications

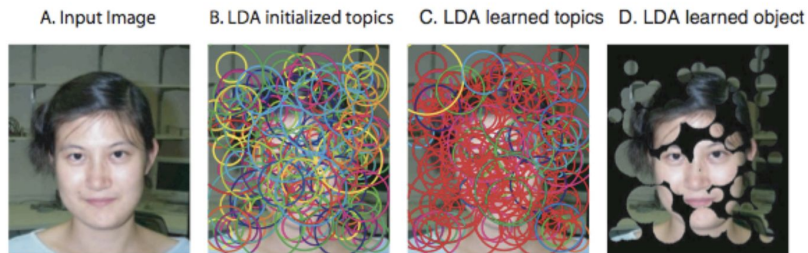# Appendix

# NLP Applications for LDA

| NLP Task | Explanation |
| --- | --- |
| **Similarity/recommendations** | Find related documents by comparing topic weight vectors (cf. NYT example above) |
| **Search** | In search engines, return more "topical" results first; SEO |
| **Word sense disambiguation** | In cases where multiple word meanings are possible, use LDA to suggest most likely meaning based on text/document topics |
| **Machine translation** | Similar to WSD - when multiple translations are possible for one word, use LDA to suggest most likely translation based on text/document topics |
| **Corpus exploration** | Find topic clusters in large corpora of literary texts, archives, etc. Popular in digital humanities research (e.g. this blog) |

# Additional applications for LDA - Images

**document** = **image**          **word** = **codeword (patch of image)**          **topic** = **object**



A. Input Image    B. LDA initialized topics    C. LDA learned topics    D. LDA learned object

In A - D above, the object learned is a person's face, but multiple objects could be learned in a single image, much like multiple topics could be learned in a single document

Before NN / deep learning era, LDA was a popular approach for image clustering, image retrieval and image relevance ranking

Sources: http://cseweb.ucsd.edu/~dhu/docs/exam09.pdf, http://pages.cs.wisc.edu/~pradheep/Clust-LDA.pdf

# Additional applications for LDA - Music

**document** = **song**     **word** = **note**     **topic** = **key**

**Key finding**: find a key for a song
**Modulation Tracking**: find a key for a segment