10. Unsupervised Techniques III: Topic Models

DS-GA 3001, Text as Data Arthur Spirling

April 3, 2018

Housekeeping

Housekeeping

Working up final homework.

Housekeeping

- Working up final homework.
- No lecture on April 16 (work on final projects)
- Speaker series: de Marneffe on "Computational Pragmatics"

(

Topic models are algorithms for discovering the main themes that pervade a large and otherwise unstructured collection of documents

Topic models are algorithms for discovering the main themes that pervade a large and otherwise unstructured collection of documents. Topic models can organize the collection according to the discovered themes.

Blei, 2012

Topic models are algorithms for discovering the main themes that pervade a large and otherwise unstructured collection of documents. Topic models can organize the collection according to the discovered themes.

Blei, 2012

Note that in social science we often use the outputs from topic models as a measurement strategy:

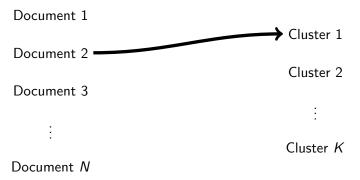
Topic models are algorithms for discovering the main themes that pervade a large and otherwise unstructured collection of documents. Topic models can organize the collection according to the discovered themes.

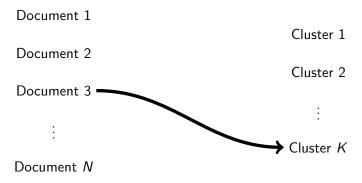
Blei, 2012

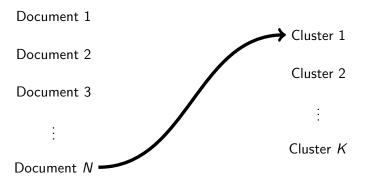
Note that in social science we often use the outputs from topic models as a measurement strategy:

"who pays more attention to education policy, conservatives or liberals?"









Document 1

Document 2

Document 3

:

Document N

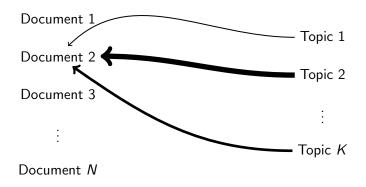
Cluster 1

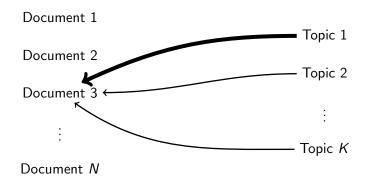
Cluster 2

:

Cluster K

| Document 1 | |
|------------|----------------|
| Document 2 | Topic 1 |
| Document 2 | Topic 2 |
| Document 3 | |
| <u>:</u> | : |
| Document N | Topic <i>K</i> |





Documents exhibit different topics,

Documents exhibit different topics, and in different proportions.

Documents exhibit different topics, and in different proportions.

E.g. a speech by the Finance minister might be 50% drawn from the trade topic, 40% from the spending topic, 9.9% from the taxation topic, 0.1% from the health topic.

Documents exhibit different topics, and in different proportions.

E.g. a speech by the Finance minister might be 50% drawn from the trade topic, 40% from the spending topic, 9.9% from the taxation topic, 0.1% from the health topic.

Think of a topic as a distribution over a fixed vocabulary.

Documents exhibit different topics, and in different proportions.

E.g. a speech by the Finance minister might be 50% drawn from the trade topic, 40% from the spending topic, 9.9% from the taxation topic, 0.1% from the health topic.

Think of a topic as a distribution over a fixed vocabulary.

E.g. the trade topic will have words like import and tariff with high probability.

Documents exhibit different topics, and in different proportions.

E.g. a speech by the Finance minister might be 50% drawn from the trade topic, 40% from the spending topic, 9.9% from the taxation topic, 0.1% from the health topic.

Think of a topic as a distribution over a fixed vocabulary.

E.g. the trade topic will have words like import and tariff with high probability.

Technically we assume the topics are generated first,

Documents exhibit different topics, and in different proportions.

E.g. a speech by the Finance minister might be 50% drawn from the trade topic, 40% from the spending topic, 9.9% from the taxation topic, 0.1% from the health topic.

Think of a topic as a distribution over a fixed vocabulary.

E.g. the trade topic will have words like import and tariff with high probability.

Technically we assume the topics are generated first, and the documents are generated second (from those topics).

Documents exhibit different topics, and in different proportions.

E.g. a speech by the Finance minister might be 50% drawn from the trade topic, 40% from the spending topic, 9.9% from the taxation topic, 0.1% from the health topic.

Think of a topic as a distribution over a fixed vocabulary.

E.g. the trade topic will have words like import and tariff with high probability.

Technically we assume the topics are generated first, and the documents are generated second (from those topics).

Now, where do the words in the documents come from?

For each document...

For each document...

• Randomly choose a distribution over topics.

For each document...

• Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.

For each document...

- Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.
- 2 Then, for every word in the document...

For each document...

- Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.
- Then, for every word in the document...
 - Randomly choose a topic from the distribution over topics from step 1.

For each document...

- Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.
- Then, for every word in the document...
 - Randomly choose a topic from the distribution over topics from step 1.
 - Randomly choose a word from the distribution over the vocabulary that the topic implies.

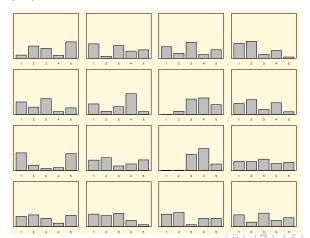
First Part

(

Randomly choose a distribution over topics.

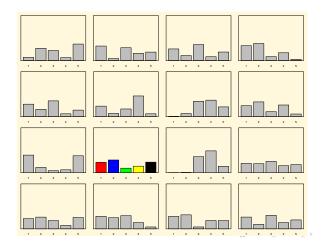
Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.

Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.

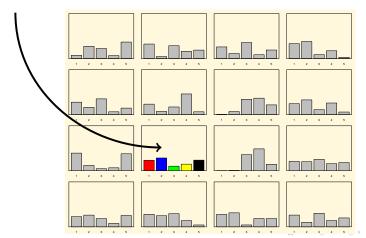


() April 3, 2018

Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.



Randomly choose a distribution over topics. That is, choose one of many multinomial distributions, each which mixes the topics in different proportions.



Then, for every word in the document...

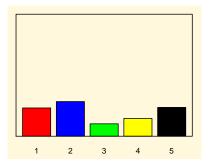
Then, for every word in the document...

• Randomly choose a topic from the distribution over topics from step 1.

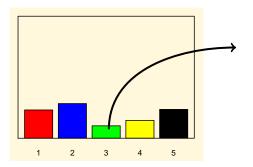
Then, for every word in the document...

- Randomly choose a topic from the distribution over topics from step 1.
- Randomly choose a word from the distribution over the vocabulary that the topic implies.

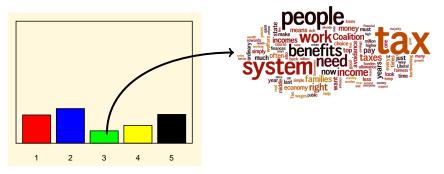
- Randomly choose a topic from the distribution over topics from step 1.
- 2 Randomly choose a word from the distribution over the vocabulary that the topic implies.



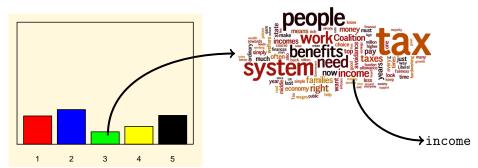
- Randomly choose a topic from the distribution over topics from step 1.
- Randomly choose a word from the distribution over the vocabulary that the topic implies.



- Randomly choose a topic from the distribution over topics from step 1.
- Randomly choose a word from the distribution over the vocabulary that the topic implies.

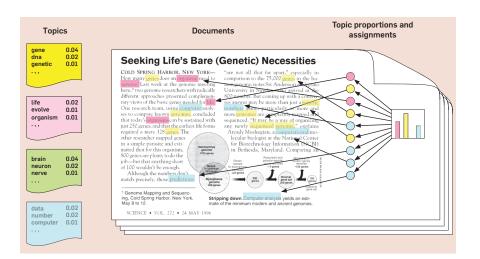


- **1** Randomly choose a topic from the distribution over topics from step 1.
- Randomly choose a word from the distribution over the vocabulary that the topic implies.



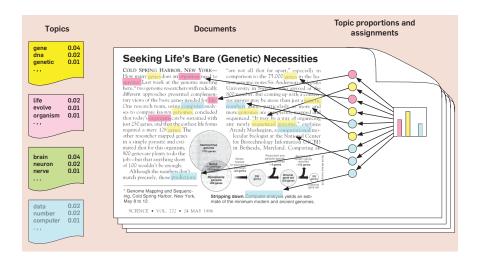
Topic Modeling a Document (Blei, 2012)

Topic Modeling a Document (Blei, 2012)



Note that all documents share same set of topics:

Topic Modeling a Document (Blei, 2012)



Note that all documents share same set of topics: but some (e.g. neuro) may be (basically) absent in a given document.

April 3, 2018

April 3, 2018

Some of our variables—the documents which contain the words—are observable.

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

We need a distribution from which to draw the per-document topic distribution. We use a Dirichlet distribution as a prior for that.

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

We need a distribution from which to draw the per-document topic distribution. We use a Dirichlet distribution as a prior for that.

And Dirichlet is used for the <u>allocation</u> of the words in the documents to different topics:

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

We need a distribution from which to draw the per-document topic distribution. We use a Dirichlet distribution as a prior for that.

And Dirichlet is used for the <u>allocation</u> of the words in the documents to different topics: it is used as a prior over the distribution of words (which define the topics).

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

We need a distribution from which to draw the per-document topic distribution. We use a Dirichlet distribution as a prior for that.

And Dirichlet is used for the <u>allocation</u> of the words in the documents to different topics: it is used as a prior over the distribution of words (which define the topics).

 \rightarrow Latent

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

We need a distribution from which to draw the per-document topic distribution. We use a Dirichlet distribution as a prior for that.

And Dirichlet is used for the <u>allocation</u> of the words in the documents to different topics: it is used as a prior over the distribution of words (which define the topics).

→ Latent Dirichlet

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

We need a distribution from which to draw the per-document topic distribution. We use a Dirichlet distribution as a prior for that.

And Dirichlet is used for the <u>allocation</u> of the words in the documents to different topics: it is used as a prior over the distribution of words (which define the topics).

→ Latent Dirichlet Allocation.

Some of our variables—the documents which contain the words—are observable. But, topic structure—topics themselves, per-document topic distributions, per-document per-word topic assignments—are latent.

We need a distribution from which to draw the per-document topic distribution. We use a Dirichlet distribution as a prior for that.

And Dirichlet is used for the <u>allocation</u> of the words in the documents to different topics: it is used as a prior over the distribution of words (which define the topics).

→ Latent Dirichlet Allocation. **LDA**



April 3, 2018

LDA is a very popular topic model:

LDA is a very popular topic model: a probabilistic procedure to generate topics. Thus, a 'generative' model.

LDA is a very popular topic model: a probabilistic procedure to generate topics. Thus, a 'generative' model.

There are *D* documents in the corpus.

LDA is a very popular topic model: a probabilistic procedure to generate topics. Thus, a 'generative' model.

There are D documents in the corpus. There are V terms in these D documents.

LDA is a very popular topic model: a probabilistic procedure to generate topics. Thus, a 'generative' model.

There are D documents in the corpus. There are V terms in these D documents. For now suppose we know the K topic distributions: there are K multinomials containing V elements each.

LDA is a very popular topic model: a probabilistic procedure to generate topics. Thus, a 'generative' model.

There are D documents in the corpus. There are V terms in these D documents. For now suppose we know the K topic distributions: there are K multinomials containing V elements each.

The multinomial distribution for the *i*th topic is denoted β_i , and $|\beta_i| = V$, meaning that the 'size' of this multinomial is equal to the number of different words in the corpus.

So, a little more formally...

So, a little more formally...

For each document...

So, a little more formally...

For each document...

lacktriangle Randomly choose a distribution over topics (multinomial of length K)

So, a little more formally...

For each document. . .

- lacktriangle Randomly choose a distribution over topics (multinomial of length K)
- 2 Then, for every word in the document...

So, a little more formally...

- lacktriangle Randomly choose a distribution over topics (multinomial of length K)
- 2 Then, for every word in the document...
 - Probabilistically draw one of the K topics from the distribution over topics from step 1. E.g. draw β_j

So, a little more formally...

For each document...

- lacktriangledown Randomly choose a distribution over topics (multinomial of length K)
- Then, for every word in the document...
 - $\begin{tabular}{ll} \bf Probabilistically draw one of the K topics from the distribution over topics from step 1. E.g. draw β_j \\ \end{tabular}$
 - **2** Probabilistically draw one of the V words from β_j

For each document...

① Draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$, where $\text{Dir}(\cdot)$ is a draw from a symmetric (uniform) Dirichlet distribution with scaling parameter, α

- Draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$, where $\text{Dir}(\cdot)$ is a draw from a symmetric (uniform) Dirichlet distribution with scaling parameter, α
- 2 Then, for every word in the document...

For each document...

- **①** Draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$, where $\text{Dir}(\cdot)$ is a draw from a symmetric (uniform) Dirichlet distribution with scaling parameter, α
- Then, for every word in the document...
 - Draw a specific topic $z_{d,n} \sim \operatorname{multi}(\theta_d)$ where $\operatorname{multi}(\cdot)$ is a multinomial. Here $z_{d,n}$ is the topic assignment for the word in the nth position of the dth document. E.g. word in position 2 in document 5 is from Topic 6.

- Draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$, where $\text{Dir}(\cdot)$ is a draw from a symmetric (uniform) Dirichlet distribution with scaling parameter, α
- Then, for every word in the document...
 - **1** Draw a specific topic $z_{d,n} \sim \operatorname{multi}(\theta_d)$ where $\operatorname{multi}(\cdot)$ is a multinomial. Here $z_{d,n}$ is the topic assignment for the word in the nth position of the dth document. E.g. word in position 2 in document 5 is from Topic 6. BTW, the multinomial has only one trial.

- **1** Draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$, where $\text{Dir}(\cdot)$ is a draw from a symmetric (uniform) Dirichlet distribution with scaling parameter, α
- Then, for every word in the document...
 - Draw a specific topic $z_{d,n} \sim \operatorname{multi}(\theta_d)$ where $\operatorname{multi}(\cdot)$ is a multinomial. Here $z_{d,n}$ is the topic assignment for the word in the nth position of the dth document. E.g. word in position 2 in document 5 is from Topic 6. BTW, the multinomial has only one trial.
 - ② Draw a word $w_{d,n} \sim \beta_{z_{d,n}}$. Here, $w_{d,n}$ is the word in the *n*th position of the *d*th document and it is being drawn from topic $\beta_{z_{d,n}}$.

- **1** Draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$, where $\text{Dir}(\cdot)$ is a draw from a symmetric (uniform) Dirichlet distribution with scaling parameter, α
- Then, for every word in the document...
 - Draw a specific topic $z_{d,n} \sim \operatorname{multi}(\theta_d)$ where $\operatorname{multi}(\cdot)$ is a multinomial. Here $z_{d,n}$ is the topic assignment for the word in the nth position of the dth document. E.g. word in position 2 in document 5 is from Topic 6. BTW, the multinomial has only one trial.
 - ② Draw a word $w_{d,n} \sim \beta_{z_{d,n}}$. Here, $w_{d,n}$ is the word in the *n*th position of the *d*th document and it is being drawn from topic $\beta_{z_{d,n}}$. E.g. word in position 2 in document 5 is from Topic 6 and turns out to be 'income' in this particular case.

The Dirichlet distribution is a conjugate prior for the multinomial ('categorical' if you only have one trial) distribution.

The Dirichlet distribution is a conjugate prior for the multinomial ('categorical' if you only have one trial) distribution. Makes certain calculations easier.

The Dirichlet distribution is a conjugate prior for the multinomial ('categorical' if you only have one trial) distribution. Makes certain calculations easier.

It is parameterized by a vector of positive real numbers α . In principle, one can have $\alpha_1, \ldots, \alpha_k$ be different concentration parameters, but LDA uses special symmetric Dirichlet where all the values of α are the same.

The Dirichlet distribution is a conjugate prior for the multinomial ('categorical' if you only have one trial) distribution. Makes certain calculations easier.

It is parameterized by a vector of positive real numbers α . In principle, one can have $\alpha_1, \ldots, \alpha_k$ be different concentration parameters, but LDA uses special symmetric Dirichlet where all the values of α are the same.

Larger values of α (assuming we are in symmetric case) mean we think (a priori) that documents are generally an even mix of the topics.

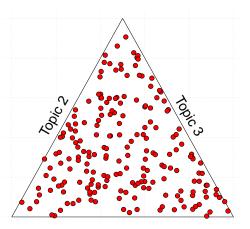
The Dirichlet distribution is a conjugate prior for the multinomial ('categorical' if you only have one trial) distribution. Makes certain calculations easier.

It is parameterized by a vector of positive real numbers α . In principle, one can have $\alpha_1, \ldots, \alpha_k$ be different concentration parameters, but LDA uses special symmetric Dirichlet where all the values of α are the same.

Larger values of α (assuming we are in symmetric case) mean we think (a priori) that documents are generally an even mix of the topics. If α is small (less than 1) we think a given document is generally from one or a few topics.

Example of Dirichlet

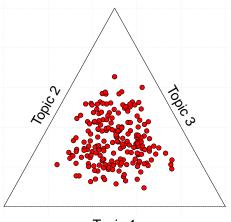
200 documents, 3 topics, $\alpha=1$ (uniform)



Topic 1

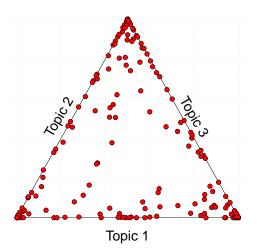
Example of Dirichlet

200 documents, 3 topics, $\alpha = 5$



Example of Dirichlet

200 documents, 3 topics, $\alpha = 0.2$



We also use a symmetric Dirichlet prior on the per topic word distributions.

We also use a symmetric Dirichlet prior on the per topic word distributions. That is, the prior on the β_i s.

We also use a symmetric Dirichlet prior on the per topic word distributions. That is, the prior on the β_i s.

 \rightarrow A high concentration parameter means each topic is a mixture of most of the words.

We also use a symmetric Dirichlet prior on the per topic word distributions. That is, the prior on the β_i s.

 \rightarrow A high concentration parameter means each topic is a mixture of most of the words. A low concentration parameter means each topic is a mixture of a few of the words.

We also use a symmetric Dirichlet prior on the per topic word distributions. That is, the prior on the β_i s.

 \rightarrow A high concentration parameter means each topic is a mixture of most of the words. A low concentration parameter means each topic is a mixture of a few of the words.

In practice, one can estimate the concentration parameters, or simple set them at suggested values.

We also use a symmetric Dirichlet prior on the per topic word distributions. That is, the prior on the β_i s.

 \rightarrow A high concentration parameter means each topic is a mixture of most of the words. A low concentration parameter means each topic is a mixture of a few of the words.

In practice, one can estimate the concentration parameters, or simple set them at suggested values.

We want topic models to be similar as we increase number of topics. Can use asymmetric priors for per-document topic distributions (the θ s).

We also use a symmetric Dirichlet prior on the per topic word distributions. That is, the prior on the β_i s.

 \rightarrow A high concentration parameter means each topic is a mixture of most of the words. A low concentration parameter means each topic is a mixture of a few of the words.

In practice, one can estimate the concentration parameters, or simple set them at suggested values.

We want topic models to be similar as we increase number of topics. Can use asymmetric priors for per-document topic distributions (the θ s). Asymmetric priors on per-topic word distributions don't do much.

We also use a symmetric Dirichlet prior on the per topic word distributions. That is, the prior on the β_i s.

 \rightarrow A high concentration parameter means each topic is a mixture of most of the words. A low concentration parameter means each topic is a mixture of a few of the words.

In practice, one can estimate the concentration parameters, or simple set them at suggested values.

We want topic models to be similar as we increase number of topics. Can use asymmetric priors for per-document topic distributions (the θ s). Asymmetric priors on per-topic word distributions don't do much. Wallach et al "Rethinking LDA: Why Priors Matter"

We observe $w_{d,n}$.

We observe $w_{d,n}$. And there are N words in a given document.

We observe $w_{d,n}$. And there are N words in a given document. And D documents in a corpus.

We observe $w_{d,n}$. And there are N words in a given document. And D documents in a corpus.

The $w_{d,n}$ we see depends on $z_{d,n}$, the topic assignment for that word ("this word will be from topic 4").

We observe $w_{d,n}$. And there are N words in a given document. And D documents in a corpus.

The $w_{d,n}$ we see depends on $z_{d,n}$, the topic assignment for that word ("this word will be from topic 4").

The $z_{d,n}$ depends on θ_d , the topic mix for a given document d in which the words sit.

We observe $w_{d,n}$. And there are N words in a given document. And D documents in a corpus.

The $w_{d,n}$ we see depends on $z_{d,n}$, the topic assignment for that word ("this word will be from topic 4").

The $z_{d,n}$ depends on θ_d , the topic mix for a given document d in which the words sit.

The θ_d depends on our prior for the relevant Dirichlet,

We now know that...

We observe $w_{d,n}$. And there are N words in a given document. And D documents in a corpus.

The $w_{d,n}$ we see depends on $z_{d,n}$, the topic assignment for that word ("this word will be from topic 4").

The $z_{d,n}$ depends on θ_d , the topic mix for a given document d in which the words sit.

The θ_d depends on our prior for the relevant Dirichlet, α .

We now know that...

We observe $w_{d,n}$. And there are N words in a given document. And D documents in a corpus.

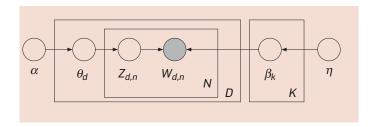
The $w_{d,n}$ we see depends on $z_{d,n}$, the topic assignment for that word ("this word will be from topic 4").

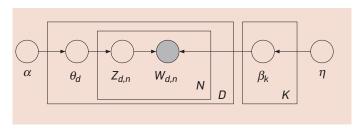
The $z_{d,n}$ depends on θ_d , the topic mix for a given document d in which the words sit.

The θ_d depends on our prior for the relevant Dirichlet, α .

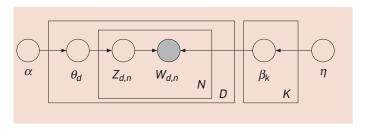
And we know that the actual value that $w_{d,n}$ takes depends on the distribution over words that the relevant topic entails, the β ("the word from topic 4 is "income" in this case")

While the β depends on the prior for the relevant Dirichlet, η

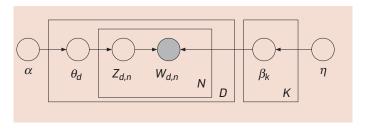




Solid nodes are observed;

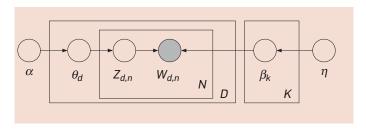


Solid nodes are observed; empty nodes are latent.



Solid nodes are observed; empty nodes are latent.

Plates imply replication.



Solid nodes are observed; empty nodes are latent.

Plates imply replication.

Note that $w_{d,n}$ depends on $z_{d,n}$ (the mix of topics for that document) and $\beta_{1:K}$ (all the topics in terms of their distributions over the words).

Bayesian Inference: Crash Course/Reminder

Recall that...

-(

Recall that...

$$\mathsf{conditional} = \frac{\mathsf{joint}}{\mathsf{marginal}}$$

Recall that...

$$conditional = \frac{joint}{marginal}$$

So,

$$Pr(A|B) = \frac{Pr(A,B)}{Pr(B)}$$

And, Bayes Theorem tell us that

$$Pr(A|B) = \frac{Pr(A) Pr(B|A)}{Pr(B)}$$

What is the probability that Republicans will win in 2020?

What is the probability that Republicans will win in 2020?

We might have a guess at this.

What is the probability that Republicans will win in 2020?

We might have a guess at this. Summarize as a prior, Pr(A).

What is the probability that Republicans will win in 2020?

We might have a guess at this. Summarize as a prior, Pr(A).

e.g. I think it is Bernoulli with p = 0.8 (so, likely, but by no means certain). You could have $\mathcal{B}(\frac{2}{3})$ instead.

What is the probability that Republicans will win in 2020?

We might have a guess at this. Summarize as a prior, Pr(A).

e.g. I think it is Bernoulli with p=0.8 (so, likely, but by no means certain). You could have $\mathcal{B}(\frac{2}{3})$ instead.

Presumably, we think the *actual* probability GOP win depends on a whole set of things:

What is the probability that Republicans will win in 2020?

We might have a guess at this. Summarize as a prior, Pr(A).

e.g. I think it is Bernoulli with p=0.8 (so, likely, but by no means certain). You could have $\mathcal{B}(\frac{2}{3})$ instead.

Presumably, we think the *actual* probability GOP win depends on a whole set of things: state of economy, success in war, etc.

What is the probability that Republicans will win in 2020?

We might have a guess at this. Summarize as a prior, Pr(A).

e.g. I think it is Bernoulli with p=0.8 (so, likely, but by no means certain). You could have $\mathcal{B}(\frac{2}{3})$ instead.

Presumably, we think the *actual* probability GOP win depends on a whole set of things: state of economy, success in war, etc. Call these data B.

What is the probability that Republicans will win in 2020?

We might have a guess at this. Summarize as a prior, Pr(A).

e.g. I think it is Bernoulli with p=0.8 (so, likely, but by no means certain). You could have $\mathcal{B}(\frac{2}{3})$ instead.

Presumably, we think the *actual* probability GOP win depends on a whole set of things: state of economy, success in war, etc. Call these data B.

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning,

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning, given observed data. It is equal to

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning, given observed data. It is equal to the product of the prior beliefs about how that probability might be distributed

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning, given observed data. It is equal to the product of the prior beliefs about how that probability might be distributed $\times \Pr(B|A)$... divided by the (unconditional) probability of the data.

April 3, 2018

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning, given observed data. It is equal to the product of the prior beliefs about how that probability might be distributed $\times \Pr(B|A)$... divided by the (unconditional) probability of the data.

We will get at Pr(B|A) via a model for the DGP.

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning, given observed data. It is equal to the product of the prior beliefs about how that probability might be distributed $\times \Pr(B|A)$... divided by the (unconditional) probability of the data.

We will get at Pr(B|A) via a model for the DGP.

We will call Pr(A|B) the posterior.

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning, given observed data. It is equal to the product of the prior beliefs about how that probability might be distributed $\times \Pr(B|A)$... divided by the (unconditional) probability of the data.

We will get at Pr(B|A) via a model for the DGP.

We will call Pr(A|B) the posterior. It will imply that some values are more plausible than others,

April 3, 2018

Interest is in
$$Pr(A|B) = \frac{Pr(A)Pr(B|A)}{Pr(B)} = Pr(win|data)$$
.

This is the estimated probability of winning, given observed data. It is equal to the product of the prior beliefs about how that probability might be distributed $\times \Pr(B|A)$... divided by the (unconditional) probability of the data.

We will get at Pr(B|A) via a model for the DGP.

We will call Pr(A|B) the posterior. It will imply that some values are more plausible than others, and we might want various features of it, like the mean or median. NB: in MLE, we would report the mode of θ

◆□ ト ◆ ② ト ◆ 夏 ト ◆ 夏 ト ◆ 夏 ◆ ② へ ○
 April 3, 2018

We have data X.

C

We have data X. We want to make an inference from it,

We have data X. We want to make an inference from it, so we assume it was produced by some (topic) model M, which has parameters θ .

We have data X. We want to make an inference from it, so we assume it was produced by some (topic) model M, which has parameters θ .

Further, assume that **X** are iid, and generated by some likelihood $p(\mathbf{X}|\boldsymbol{\theta}, M)$.

More formally

We have data X. We want to make an inference from it, so we assume it was produced by some (topic) model M, which has parameters θ .

Further, assume that **X** are iid, and generated by some likelihood $p(\mathbf{X}|\boldsymbol{\theta}, M)$.

The posterior over the parameters = likelihood of the data, conditioned on particular values for the parameters, multiplied by our prior.

More formally

We have data X. We want to make an inference from it, so we assume it was produced by some (topic) model M, which has parameters θ .

Further, assume that **X** are iid, and generated by some likelihood $p(\mathbf{X}|\boldsymbol{\theta}, M)$.

The posterior over the parameters = likelihood of the data, conditioned on particular values for the parameters, multiplied by our prior.

Then:
$$Pr(A|B) = \frac{Pr(A) Pr(B|A)}{Pr(B)}$$
.

Then:
$$Pr(A|B) = \frac{Pr(A) Pr(B|A)}{Pr(B)}$$
.

Then:
$$Pr(A|B) = \frac{Pr(A) \frac{Pr(B|A)}{Pr(B)}}{Pr(B)}$$
.

$$p(\theta|\mathbf{X}, M) = \frac{p(\mathbf{X}, \theta|M)}{\int p(\mathbf{X}, \theta|M) d\theta} = \boxed{\frac{p(\theta|M)p(\mathbf{X}|\theta, M)}{p(\mathbf{X}|M)}}$$

We refer to denominator as the 'normalizing constant' or the 'marginal likelihood'—

Then:
$$Pr(A|B) = \frac{Pr(A) \frac{Pr(B|A)}{Pr(B)}}{Pr(B)}$$
.

$$p(\theta|\mathbf{X}, M) = \frac{p(\mathbf{X}, \theta|M)}{\int p(\mathbf{X}, \theta|M) d\theta} = \boxed{\frac{p(\theta|M)p(\mathbf{X}|\theta, M)}{p(\mathbf{X}|M)}}$$

We refer to denominator as the 'normalizing constant' or the 'marginal likelihood'—because it's the likelihood with the model parameters integrated out.

Then:
$$Pr(A|B) = \frac{Pr(A) Pr(B|A)}{Pr(B)}$$
.

$$p(\theta|\mathbf{X}, M) = \frac{p(\mathbf{X}, \theta|M)}{\int p(\mathbf{X}, \theta|M) d\theta} = \boxed{\frac{p(\theta|M)p(\mathbf{X}|\theta, M)}{p(\mathbf{X}|M)}}$$

We refer to denominator as the 'normalizing constant' or the 'marginal likelihood'—because it's the likelihood with the model parameters integrated out.

NB: sometimes called the evidence in the Bayesian context, and is integral of numerator over support of θ (weighted by how plausible each value is)

April 3, 2018

(

1 Predicting values of new data points X_{new} given the observed data.

Predicting values of new data points X_{new} given the observed data.
We have to average over the posterior:
(Note: 100 to 100

 $p(\mathbf{X}_{\mathsf{new}}|\mathbf{X}, M) = \int p(\mathbf{X}_{\mathsf{new}}|\boldsymbol{\theta}, M)p(\boldsymbol{\theta}|\mathbf{X}, M)d\boldsymbol{\theta}$

Predicting values of new data points X_{new} given the observed data. We have to average over the posterior:

$$p(\mathbf{X}_{new}|\mathbf{X}, M) = \int p(\mathbf{X}_{new}|\boldsymbol{\theta}, M)p(\boldsymbol{\theta}|\mathbf{X}, M)d\boldsymbol{\theta}$$

But, this integral is often intractable because posterior, $p(\theta|\mathbf{X}, M)$ is very high dimensional or has awkward form.

Predicting values of new data points \mathbf{X}_{new} given the observed data. We have to average over the posterior: $p(\mathbf{X}_{\text{new}}|\mathbf{X},M) = \int p(\mathbf{X}_{\text{new}}|\boldsymbol{\theta},M)p(\boldsymbol{\theta}|\mathbf{X},M)d\boldsymbol{\theta}$ But, this integral is often intractable because posterior, $p(\boldsymbol{\theta}|\mathbf{X},M)$ is very high dimensional or has awkward form.

Marginalization:

Predicting values of new data points X_{new} given the observed data. We have to average over the posterior: p(X_{new}|X, M) = ∫ p(X_{new}|θ, M)p(θ|X, M)dθ
But, this integral is often intractable because posterior, p(θ|X, M) is very high dimensional or has awkward form.

@ Marginalization: may have θ and another 'nuisance' parameter \mathbf{z} ,

Predicting values of new data points \mathbf{X}_{new} given the observed data. We have to average over the posterior: $p(\mathbf{X}_{\text{new}}|\mathbf{X},M) = \int p(\mathbf{X}_{\text{new}}|\boldsymbol{\theta},M)p(\boldsymbol{\theta}|\mathbf{X},M)d\boldsymbol{\theta}$ But, this integral is often intractable because posterior, $p(\boldsymbol{\theta}|\mathbf{X},M)$ is very high dimensional or has awkward form.

Marginalization: may have θ and another 'nuisance' parameter \mathbf{z} , and we want $p(\theta|\mathbf{X}, M) = \int p(\theta, \mathbf{z}|\mathbf{X}, M) d\mathbf{z}$.

Predicting values of new data points \mathbf{X}_{new} given the observed data. We have to average over the posterior: $p(\mathbf{X}_{\text{new}}|\mathbf{X},M) = \int p(\mathbf{X}_{\text{new}}|\boldsymbol{\theta},M)p(\boldsymbol{\theta}|\mathbf{X},M)d\boldsymbol{\theta}$ But, this integral is often intractable because posterior, $p(\boldsymbol{\theta}|\mathbf{X},M)$ is very high dimensional or has awkward form.

Marginalization: may have θ and another 'nuisance' parameter \mathbf{z} , and we want $p(\theta|\mathbf{X},M) = \int p(\theta,\mathbf{z}|\mathbf{X},M)d\mathbf{z}$. Which may require intractable integration.

April 3, 2018

Predicting values of new data points \mathbf{X}_{new} given the observed data. We have to average over the posterior: $p(\mathbf{X}_{\text{new}}|\mathbf{X},M) = \int p(\mathbf{X}_{\text{new}}|\boldsymbol{\theta},M)p(\boldsymbol{\theta}|\mathbf{X},M)d\boldsymbol{\theta}$ But, this integral is often intractable because posterior, $p(\boldsymbol{\theta}|\mathbf{X},M)$ is very high dimensional or has awkward form.

- **Marginalization:** may have θ and another 'nuisance' parameter \mathbf{z} , and we want $p(\theta|\mathbf{X},M) = \int p(\theta,\mathbf{z}|\mathbf{X},M)d\mathbf{z}$. Which may require intractable integration.
- Model selection:

April 3, 2018

- Predicting values of new data points \mathbf{X}_{new} given the observed data. We have to average over the posterior: $p(\mathbf{X}_{\text{new}}|\mathbf{X},M) = \int p(\mathbf{X}_{\text{new}}|\boldsymbol{\theta},M)p(\boldsymbol{\theta}|\mathbf{X},M)d\boldsymbol{\theta}$ But, this integral is often intractable because posterior, $p(\boldsymbol{\theta}|\mathbf{X},M)$ is very high dimensional or has awkward form.
- **Marginalization:** may have θ and another 'nuisance' parameter \mathbf{z} , and we want $p(\theta|\mathbf{X},M) = \int p(\theta,\mathbf{z}|\mathbf{X},M)d\mathbf{z}$. Which may require intractable integration.
- Model selection: might want to compare models, and see which is most plausible.

April 3, 2018

- Predicting values of new data points \mathbf{X}_{new} given the observed data. We have to average over the posterior: $p(\mathbf{X}_{\text{new}}|\mathbf{X},M) = \int p(\mathbf{X}_{\text{new}}|\boldsymbol{\theta},M)p(\boldsymbol{\theta}|\mathbf{X},M)d\boldsymbol{\theta}$ But, this integral is often intractable because posterior, $p(\boldsymbol{\theta}|\mathbf{X},M)$ is very high dimensional or has awkward form.
- **Marginalization:** may have θ and another 'nuisance' parameter \mathbf{z} , and we want $p(\theta|\mathbf{X},M) = \int p(\theta,\mathbf{z}|\mathbf{X},M)d\mathbf{z}$. Which may require intractable integration.
- **Model selection**: might want to compare models, and see which is most plausible. But that requires calculating $\int p(\mathbf{X}, \theta | M) d\theta$

0

- Predicting values of new data points \mathbf{X}_{new} given the observed data. We have to average over the posterior: $p(\mathbf{X}_{\text{new}}|\mathbf{X},M) = \int p(\mathbf{X}_{\text{new}}|\boldsymbol{\theta},M)p(\boldsymbol{\theta}|\mathbf{X},M)d\boldsymbol{\theta}$ But, this integral is often intractable because posterior, $p(\boldsymbol{\theta}|\mathbf{X},M)$ is very high dimensional or has awkward form.
- **Marginalization:** may have θ and another 'nuisance' parameter \mathbf{z} , and we want $p(\theta|\mathbf{X},M) = \int p(\theta,\mathbf{z}|\mathbf{X},M)d\mathbf{z}$. Which may require intractable integration.
- **Model selection**: might want to compare models, and see which is most plausible. But that requires calculating $\int p(\mathbf{X}, \boldsymbol{\theta}|M) d\boldsymbol{\theta}$ ('evidence') which may be intractable.

April 3, 2018

April 3, 2018

Bayesian approaches require integration.

Bayesian approaches require integration. But those integrals may be difficult.

Bayesian approaches require integration. But those integrals may be difficult.

Instead, we can approximate. Two main ideas:

Bayesian approaches require integration. But those integrals may be difficult.

Instead, we can approximate. Two main ideas:

Stochastic: we can use Monte Carlo simulation methods to sample from the posterior distribution, e.g. $p(\theta|\mathbf{X}, M)$.

Bayesian approaches require integration. But those integrals may be difficult.

Instead, we can approximate. Two main ideas:

Stochastic: we can use Monte Carlo simulation methods to sample from the posterior distribution, e.g. $p(\theta|\mathbf{X}, M)$. Markov chain Monte Carlo is a way to do this,

Bayesian approaches require integration. But those integrals may be difficult.

Instead, we can approximate. Two main ideas:

Stochastic: we can use Monte Carlo simulation methods to sample from the posterior distribution, e.g. $p(\theta|\mathbf{X}, M)$. Markov chain Monte Carlo is a way to do this, with the Gibbs Sampler a specific example.

Bayesian approaches require integration. But those integrals may be difficult.

Instead, we can approximate. Two main ideas:

Stochastic: we can use Monte Carlo simulation methods to sample from the posterior distribution, e.g. $p(\theta|\mathbf{X}, M)$. Markov chain Monte Carlo is a way to do this, with the Gibbs Sampler a specific example.

Deterministic: most notably for topic models, variational inference.

Bayesian approaches require integration. But those integrals may be difficult.

Instead, we can approximate. Two main ideas:

Stochastic: we can use Monte Carlo simulation methods to sample from the posterior distribution, e.g. $p(\theta|\mathbf{X}, M)$. Markov chain Monte Carlo is a way to do this, with the Gibbs Sampler a specific example.

Deterministic: most notably for topic models, variational inference. Idea is to write down posterior as product of well-known distributions.

Bayesian approaches require integration. But those integrals may be difficult.

Instead, we can approximate. Two main ideas:

Stochastic: we can use Monte Carlo simulation methods to sample from the posterior distribution, e.g. $p(\theta|\mathbf{X}, M)$. Markov chain Monte Carlo is a way to do this, with the Gibbs Sampler a specific example.

Deterministic: most notably for topic models, variational inference. Idea is to write down posterior as product of well-known distributions. This will approximate the true posterior (use KL divergence to get as close as possible to it).

Bayesian computation is difficult and intensive compared to, say, MLE.

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy:

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense.

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena.

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process.

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data.

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty.

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty. Data overwhelms prior.

Bayesian computation is difficult and intensive compared to, say, MLE. So why do it?

1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with p-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty. Data overwhelms prior.

2 Practicality:

- 1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty. Data overwhelms prior.
- 2 Practicality: for many problems, MLE doesn't work well.

- 1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty. Data overwhelms prior.
- 2 Practicality: for many problems, MLE doesn't work well. Perhaps because there are lots of parameters, but not much data:

- 1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty. Data overwhelms prior.
- 2 Practicality: for many problems, MLE doesn't work well. Perhaps because there are lots of parameters, but not much data: priors can help here.

- 1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty. Data overwhelms prior.
- 2 Practicality: for many problems, MLE doesn't work well. Perhaps because there are lots of parameters, but not much data: priors can help here. Or perhaps because the model, like multinomial probit, involves evaluating something complicated, analytically:

- 1 Philosophy: Bayesian (v frequentist) approach to probability and inference simply makes more sense. We do away with *p*-values and notions of repeating (finite) phenomena. We have priors from previous work and can include them explicitly as part of scientific process. Automatic handling of missing data. More sensible handling of uncertainty. Data overwhelms prior.
- 2 Practicality: for many problems, MLE doesn't work well. Perhaps because there are lots of parameters, but not much data: priors can help here. Or perhaps because the model, like multinomial probit, involves evaluating something complicated, analytically: Bayesian methods can arbitrarily approximate the integral.

Crash course complete: back to LDA

Ultimately,

C

Ultimately, we will use the observed data, the words,

Ultimately, we will use the observed data, the words, to make an inference about the latent parameters:

Ultimately, we will use the observed data, the words, to make an inference about the latent parameters: the β s, the zs, the θ s.

Ultimately, we will use the observed data, the words, to make an inference about the latent parameters: the β s, the zs, the θ s. That will be a conditional probability.

Ultimately, we will use the observed data, the words, to make an inference about the latent parameters: the β s, the zs, the θ s. That will be a conditional probability.

We start with the joint distribution implied by the problem:

Ultimately, we will use the observed data, the words, to make an inference about the latent parameters: the β s, the zs, the θ s. That will be a conditional probability.

We start with the joint distribution implied by the problem:

$$p(\beta_{1:K},\theta_{1:D},z_{1:D},w_{1:D}) =$$

Ultimately, we will use the observed data, the words, to make an inference about the latent parameters: the β s, the zs, the θ s. That will be a conditional probability.

We start with the joint distribution implied by the problem:

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D}) =$$

$$\prod_{K}^{i=1} p(\beta_i) \prod_{D}^{d=1} p(\theta_d) \left(\prod_{N}^{n=1} p(z_{d,n}|\theta_d) p(w_{d,n}|\beta_{1:K}, z_{d,n}) \right)$$

Generally we want

$$p(\boldsymbol{\theta}|\mathbf{X}, M) = \frac{p(\mathbf{X}, \boldsymbol{\theta}|M)}{\int p(\mathbf{X}, \boldsymbol{\theta}|M) d\boldsymbol{\theta}} = \frac{p(\boldsymbol{\theta}|M)p(\mathbf{X}|\boldsymbol{\theta}, M)}{p(\mathbf{X}|M)}$$

Generally we want

$$p(\theta|\mathbf{X}, M) = \frac{p(\mathbf{X}, \theta|M)}{\int p(\mathbf{X}, \theta|M)d\theta} = \boxed{\frac{p(\theta|M)p(\mathbf{X}|\theta, M)}{p(\mathbf{X}|M)}}$$

Here (Blei, 2012), that is

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}|w_{1:D}) = \frac{p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})}{p(w_{1:D})}$$

Generally we want

$$p(\theta|\mathbf{X}, M) = \frac{p(\mathbf{X}, \theta|M)}{\int p(\mathbf{X}, \theta|M)d\theta} = \frac{p(\theta|M)p(\mathbf{X}|\theta, M)}{p(\mathbf{X}|M)}$$

Here (Blei, 2012), that is

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}|w_{1:D}) = \frac{p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})}{p(w_{1:D})}$$

Can get 'evidence' (denominator) by summing joint distribution over every possible topic structure:

Generally we want

$$p(\theta|\mathbf{X}, M) = \frac{p(\mathbf{X}, \theta|M)}{\int p(\mathbf{X}, \theta|M)d\theta} = \boxed{\frac{p(\theta|M)p(\mathbf{X}|\theta, M)}{p(\mathbf{X}|M)}}$$

Here (Blei, 2012), that is

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}|w_{1:D}) = \frac{p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})}{p(w_{1:D})}$$

Can get 'evidence' (denominator) by summing joint distribution over every possible topic structure: every possible way of assigning each word to a topic.

Generally we want

$$p(\theta|\mathbf{X}, M) = \frac{p(\mathbf{X}, \theta|M)}{\int p(\mathbf{X}, \theta|M)d\theta} = \boxed{\frac{p(\theta|M)p(\mathbf{X}|\theta, M)}{p(\mathbf{X}|M)}}$$

Here (Blei, 2012), that is

$$p(\beta_{1:K}, \theta_{1:D}, z_{1:D}|w_{1:D}) = \frac{p(\beta_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D})}{p(w_{1:D})}$$

Can get 'evidence' (denominator) by summing joint distribution over every possible topic structure: every possible way of assigning each word to a topic. But this is impossible, so simulate/approximate.

For a user-selected k, a typical implementation of LDA will return...

For a user-selected k, a typical implementation of LDA will return...

The word distribution for each topic.

For a user-selected k, a typical implementation of LDA will return...

The word distribution for each topic.

The topic distribution for each document.

For a user-selected k, a typical implementation of LDA will return...

The word distribution for each topic.

The topic distribution for each document.

Some implementations allow you to estimate e.g. α , in which case this is also returned.

For a user-selected k, a typical implementation of LDA will return...

The word distribution for each topic.

The topic distribution for each document.

Some implementations allow you to estimate e.g. α , in which case this is also returned. And perhaps some kind of fit statistic(s).

69 UK manifestos.

69 UK manifestos. Some preprocessing.

69 UK manifestos. Some preprocessing. Used topicmodels to fit five topics.

69 UK manifestos. Some preprocessing. Used topicmodels to fit five topics. Has Gibbs sampling and variational options.

69 UK manifestos. Some preprocessing. Used topicmodels to fit five topics. Has Gibbs sampling and variational options.

The (some selected) word distributions for each topic.

69 UK manifestos. Some preprocessing. Used topicmodels to fit five topics. Has Gibbs sampling and variational options.

The (some selected) word distributions for each topic. Sum down the columns is one.

69 UK manifestos. Some preprocessing. Used topicmodels to fit five topics. Has Gibbs sampling and variational options.

The (some selected) word distributions for each topic. Sum down the columns is one.

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|--------------|---------|---------|---------|---------|---------|
| conservative | 0.00188 | 0.00088 | 0.00185 | 0.00221 | 0.00168 |
| party | 0.00145 | 0.00067 | 0.00066 | 0.00577 | 0.00093 |
| general | 0.00073 | 0.00033 | 0.00018 | 0.00192 | 0.00040 |
| election | 0.00079 | 0.00053 | 0.00022 | 0.00235 | 0.00076 |
| manifesto | 0.00059 | 0.00078 | 0.00032 | 0.00099 | 0.00048 |
| : | : | | : | : | : |
| : | : | | : | : | : |

'Top' 6 most frequent words in each topic:

'Top' 6 most frequent words in each topic: might help interpretation (!)

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---|------------|--------------|----------|------------|------------|
| 1 | people | new | [markup] | new | must |
| 2 | local | government | people | labour | government |
| 3 | government | people | new | government | labour |
| 4 | new | continue | work | people | shall |
| 5 | tax | can | [markup] | shall | can |
| 6 | liberal | conservative | support | britain | policy |

'Top' 6 most frequent words in each topic: might help interpretation (!)

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---|------------|--------------|----------|------------|------------|
| 1 | people | new | [markup] | new | must |
| 2 | local | government | people | labour | government |
| 3 | government | people | new | government | labour |
| 4 | new | continue | work | people | shall |
| 5 | tax | can | [markup] | shall | can |
| 6 | liberal | conservative | support | britain | policy |

Up to analyst to label the topics!

'Top' 6 most frequent words in each topic: might help interpretation (!)

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---|------------|--------------|----------|------------|------------|
| 1 | people | new | [markup] | new | must |
| 2 | local | government | people | labour | government |
| 3 | government | people | new | government | labour |
| 4 | new | continue | work | people | shall |
| 5 | tax | can | [markup] | shall | can |
| 6 | liberal | conservative | support | britain | policy |

Up to analyst to label the topics!

Meaningless 'junk' topics not unusual:

'Top' 6 most frequent words in each topic: might help interpretation (!)

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---|------------|--------------|----------|------------|------------|
| 1 | people | new | [markup] | new | must |
| 2 | local | government | people | labour | government |
| 3 | government | people | new | government | labour |
| 4 | new | continue | work | people | shall |
| 5 | tax | can | [markup] | shall | can |
| 6 | liberal | conservative | support | britain | policy |

Up to analyst to label the topics!

Meaningless 'junk' topics not unusual: debate as to whether one has to interpret every topic.

The topic distribution for each document...

The topic distribution for each document...

The topic distribution for each document...

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|-------|---------|---------|---------|---------|---------|
| doc 1 | 0.00009 | 0.00009 | 0.00009 | 0.00009 | 0.99965 |
| doc 2 | 0.00011 | 0.00011 | 0.00011 | 0.00011 | 0.99954 |
| doc 3 | 0.00010 | 0.00010 | 0.00010 | 0.00010 | 0.99959 |
| doc 4 | 0.00006 | 0.00006 | 0.00006 | 0.00006 | 0.99978 |
| doc 5 | 0.00002 | 0.00002 | 0.00002 | 0.00002 | 0.99991 |
| doc 6 | 0.00019 | 0.00019 | 0.00019 | 0.00019 | 0.99924 |
| | | | | | |
| i | | : | | | : |

The topic distribution for each document...

| | Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|-------|---------|---------|---------|---------|---------|
| doc 1 | 0.00009 | 0.00009 | 0.00009 | 0.00009 | 0.99965 |
| doc 2 | 0.00011 | 0.00011 | 0.00011 | 0.00011 | 0.99954 |
| doc 3 | 0.00010 | 0.00010 | 0.00010 | 0.00010 | 0.99959 |
| doc 4 | 0.00006 | 0.00006 | 0.00006 | 0.00006 | 0.99978 |
| doc 5 | 0.00002 | 0.00002 | 0.00002 | 0.00002 | 0.99991 |
| doc 6 | 0.00019 | 0.00019 | 0.00019 | 0.00019 | 0.99924 |
| | | | | | |
| i | | : | | | : |

Texts are usually preprocessed:

Texts are usually preprocessed: stop words removed,

Texts are usually preprocessed: stop words removed, (very) rare tokens removed.

Texts are usually preprocessed: stop words removed, (very) rare tokens removed. Punctuation often removed.

Texts are usually preprocessed: stop words removed, (very) rare tokens removed. Punctuation often removed. Stemming seems less common.

Texts are usually preprocessed: stop words removed, (very) rare tokens removed. Punctuation often removed. Stemming seems less common.

In most social science examples, the number of topics, K, is not picked automatically.

Texts are usually preprocessed: stop words removed, (very) rare tokens removed. Punctuation often removed. Stemming seems less common.

In most social science examples, the number of topics, K, is not picked automatically. Analysts select various Ks and check that their results are 'robust'. But see over.

Texts are usually preprocessed: stop words removed, (very) rare tokens removed. Punctuation often removed. Stemming seems less common.

In most social science examples, the number of topics, K, is not picked automatically. Analysts select various Ks and check that their results are 'robust'. But see over.

As with all unsupervised learning,

Texts are usually preprocessed: stop words removed, (very) rare tokens removed. Punctuation often removed. Stemming seems less common.

In most social science examples, the number of topics, K, is not picked automatically. Analysts select various Ks and check that their results are 'robust'. But see over

As with all unsupervised learning, interpretation is non-trivial, and requires a lot of validation. Rant: 'just-so' stories abound. Lazy analysts conclude whatever they want.

Texts are usually preprocessed: stop words removed, (very) rare tokens removed. Punctuation often removed. Stemming seems less common.

In most social science examples, the number of topics, K, is not picked automatically. Analysts select various Ks and check that their results are 'robust'. But see over

As with all unsupervised learning, interpretation is non-trivial, and requires a lot of validation. Rant: 'just-so' stories abound. Lazy analysts conclude whatever they want.

Crudely: in social science,

Crudely: in social science, researchers fit 'enough' topics until they see what they think they should.

Crudely: in social science, researchers fit 'enough' topics until they see what they think they should. E.g. a certain topic—like finance suddenly peels off—so stop there.

Crudely: in social science, researchers fit 'enough' topics until they see what they think they should. E.g. a certain topic—like finance suddenly peels off—so stop there.

 \rightarrow Check findings are robust in the neighborhood: if best model has k=35, check k=30-40 yields similar inferences.

Crudely: in social science, researchers fit 'enough' topics until they see what they think they should. E.g. a certain topic—like finance suddenly peels off—so stop there.

 \rightarrow Check findings are robust in the neighborhood: if best model has k=35, check k=30-40 yields similar inferences.

NB: social scientists typically fit far fewer topics than CS, even to same data.

Crudely: in social science, researchers fit 'enough' topics until they see what they think they should. E.g. a certain topic—like finance suddenly peels off—so stop there.

 \rightarrow Check findings are robust in the neighborhood: if best model has k=35, check k=30-40 yields similar inferences.

NB: social scientists typically fit far fewer topics than CS, even to same data.

Picking *k*, continued...

Picking *k*, continued...

CS: split into training and test sets.

Picking *k*, continued...

CS: split into training and test sets. In the training set,

CS: split into training and test sets. In the training set,

 \bigcirc pick some value of k and fit a topic model.

CS: split into training and test sets. In the training set,

- $oldsymbol{0}$ pick some value of k and fit a topic model.
- 2 record value of α (hyperparameter on document specific topic distributions) and word distributions for the topics (the β s)

CS: split into training and test sets. In the training set,

- lacksquare pick some value of k and fit a topic model.
- 2 record value of α (hyperparameter on document specific topic distributions) and word distributions for the topics (the β s)

We'll write the β s as β , then we want

$$\mathcal{L}(\mathbf{w}) = \log p(\mathbf{w}|\boldsymbol{\beta}, \alpha) = \sum_{d} \log p(w_{d}|\boldsymbol{\beta}, \alpha)$$

CS: split into training and test sets. In the training set,

- lacksquare pick some value of k and fit a topic model.
- 2 record value of α (hyperparameter on document specific topic distributions) and word distributions for the topics (the β s)

We'll write the β s as β , then we want

$$\mathcal{L}(\mathbf{w}) = \log p(\mathbf{w}|\boldsymbol{\beta}, \alpha) = \sum_{d} \log p(w_{d}|\boldsymbol{\beta}, \alpha)$$

where w are the words in the test set.

CS: split into training and test sets. In the training set,

- lacksquare pick some value of k and fit a topic model.
- 2 record value of α (hyperparameter on document specific topic distributions) and word distributions for the topics (the β s)

We'll write the β s as β , then we want

$$\mathcal{L}(\mathbf{w}) = \log p(\mathbf{w}|\boldsymbol{\beta}, \alpha) = \sum_{d} \log p(w_{d}|\boldsymbol{\beta}, \alpha)$$

where \boldsymbol{w} are the words in the test set. Higher $\boldsymbol{\mathcal{L}}$ implies better model.

CS: split into training and test sets. In the training set,

- lacktriangledown pick some value of k and fit a topic model.
- 2 record value of α (hyperparameter on document specific topic distributions) and word distributions for the topics (the β s)

We'll write the β s as β , then we want

$$\mathcal{L}(\mathbf{w}) = \log p(\mathbf{w}|\boldsymbol{\beta}, \alpha) = \sum_{d} \log p(w_{d}|\boldsymbol{\beta}, \alpha)$$

where \mathbf{w} are the words in the test set. Higher \mathcal{L} implies better model. Intuition is to calculate likelihood of seeing the test words, given what we know produced the training set.

CS: split into training and test sets. In the training set,

- lacktriangledown pick some value of k and fit a topic model.
- 2 record value of α (hyperparameter on document specific topic distributions) and word distributions for the topics (the β s)

We'll write the β s as β , then we want

$$\mathcal{L}(\mathbf{w}) = \log p(\mathbf{w}|\boldsymbol{\beta}, \alpha) = \sum_{d} \log p(w_{d}|\boldsymbol{\beta}, \alpha)$$

where \mathbf{w} are the words in the test set. Higher \mathcal{L} implies better model. Intuition is to calculate likelihood of seeing the test words, given what we know produced the training set.

Do this for all k.

CS: split into training and test sets. In the training set,

- lacktriangledown pick some value of k and fit a topic model.
- 2 record value of α (hyperparameter on document specific topic distributions) and word distributions for the topics (the β s)

We'll write the β s as β , then we want

$$\mathcal{L}(\mathbf{w}) = \log p(\mathbf{w}|\boldsymbol{\beta}, \alpha) = \sum_{d} \log p(w_{d}|\boldsymbol{\beta}, \alpha)$$

where \mathbf{w} are the words in the test set. Higher \mathcal{L} implies better model. Intuition is to calculate likelihood of seeing the test words, given what we know produced the training set.

Do this for all k.

Perplexity is popular option

Perplexity is popular option

$$\mathsf{perplexity} = \mathsf{exp}\left(-\frac{\mathcal{L}(\mathbf{w})}{\mathsf{count}\ \mathsf{of}\ \mathsf{tokens}}\right),$$

Perplexity is popular option

$$\mathsf{perplexity} = \mathsf{exp}\left(-\frac{\mathcal{L}(\mathbf{w})}{\mathsf{count}\ \mathsf{of}\ \mathsf{tokens}}\right),$$

where lower is better.

Perplexity is popular option

$$perplexity = exp\left(-\frac{\mathcal{L}(\mathbf{w})}{count of tokens}\right),$$

where lower is better.

In general, $\mathcal{L}(\mathbf{w})$ is intractable,

Perplexity is popular option

$$\mathsf{perplexity} = \mathsf{exp}\left(-\frac{\mathcal{L}(\mathbf{w})}{\mathsf{count}\ \mathsf{of}\ \mathsf{tokens}}\right),$$

where lower is better.

In general, $\mathcal{L}(\mathbf{w})$ is intractable, but there are ways to approximate it.

Perplexity is popular option

$$\mathsf{perplexity} = \mathsf{exp}\left(-\frac{\mathcal{L}(\mathbf{w})}{\mathsf{count}\ \mathsf{of}\ \mathsf{tokens}}\right),$$

where lower is better.

In general, $\mathcal{L}(\mathbf{w})$ is intractable, but there are ways to approximate it.

But:

Perplexity is popular option

$$\mathsf{perplexity} = \mathsf{exp}\left(-\frac{\mathcal{L}(\mathbf{w})}{\mathsf{count}\ \mathsf{of}\ \mathsf{tokens}}\right),$$

where lower is better.

In general, $\mathcal{L}(\mathbf{w})$ is intractable, but there are ways to approximate it.

But: the topic models that hold-out calculations suggest are optimal and not much liked by humans!

Perplexity is popular option

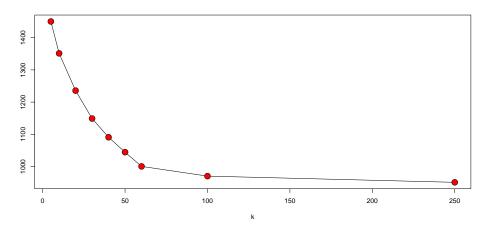
$$\mathsf{perplexity} = \mathsf{exp}\left(-\frac{\mathcal{L}(\mathbf{w})}{\mathsf{count}\ \mathsf{of}\ \mathsf{tokens}}\right),$$

where lower is better.

In general, $\mathcal{L}(\mathbf{w})$ is intractable, but there are ways to approximate it.

But: the topic models that hold-out calculations suggest are optimal and not much liked by humans! "Reading Tea Leaves: How Humans Interpret Topic Models" by Chang et al.

Perplexity Likes a Lot of Topics (manifestos)





-0



Japan is a curious IR case:



Japan is a curious IR case: wealthy post-war not very interested in foreign policy.



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area.



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?

Rise of China?



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?

1 Rise of China? Need to focus on security.



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?

• Rise of China? Need to focus on security.

VS.

2 Change in Electoral System?



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?

• Rise of China? Need to focus on security.

VS.

Change in Electoral System? Moved from promising pork to having to deliver policy as part of Westminster-style polity.



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?

• Rise of China? Need to focus on security.

VS.

Change in Electoral System? Moved from promising pork to having to deliver policy as part of Westminster-style polity.

To decide, we need data source that covers all lower house legislators

- (,



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?

- Rise of China? Need to focus on security.
- VS.
- Change in Electoral System? Moved from promising pork to having to deliver policy as part of Westminster-style polity.

To decide, we need data source that covers all lower house legislators where they set out their policy priorities over time.



Japan is a curious IR case: wealthy post-war not very interested in foreign policy. Recent times have seen a (re-)emergence in this area. Why?

- Rise of China? Need to focus on security.
- VS.
- Change in Electoral System? Moved from promising pork to having to deliver policy as part of Westminster-style polity.

To decide, we need data source that covers all lower house legislators where they set out their policy priorities over time. See if/when they shift priorities.

(

-(





7,497.

-()



7,497. 1986-2009.



7,497. 1986-2009. Standardized form.



7,497. 1986-2009. Standardized form.

"... instructed to write whatever they want in the form and return it before 5 PM of the first day of the campaign. At least two days before the election, local electoral commissions are required to distribute the forms of all candidates running in the district to all registered voters"

April 3, 2018



7,497. 1986-2009. Standardized form.

"... instructed to write whatever they want in the form and return it before 5 PM of the first day of the campaign. At least two days before the election, local electoral commissions are required to distribute the forms of all candidates running in the district to all registered voters"

Manifestos were hand transcribed from microfilm.

◆□ → ◆□ → ◆三 → □ → ○○ ○



7,497. 1986-2009. Standardized form.

"... instructed to write whatever they want in the form and return it before 5 PM of the first day of the campaign. At least two days before the election, local electoral commissions are required to distribute the forms of all candidates running in the district to all registered voters"

Manifestos were hand transcribed from microfilm. Japanese install of Windows/R used to fit LDA.

April 3, 2018

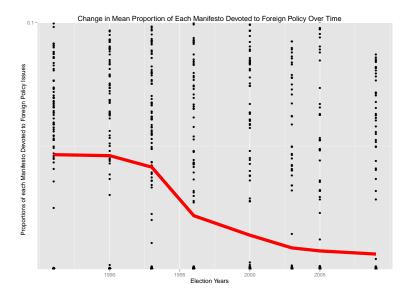
Topic Distribution over Words

Topic Distribution over Words

| Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 | Topic 6 |
|---------|---------|---------|---------|---------------|---------------|
| 1改革 | 年金 | 推進 | X | 政治 | 日本 |
| 2 郵政 | 円 | 整備 | 政策 | 改革 | 1 |
| 3 民営 | 廃止 | 図る | 地域 | 国民 | 外交 |
| 4 小泉 | 改革 | つとめる | まち | 企業 | 国家 |
| 5 構造 | 兆 | 社会 | 鹿児島 | 自民党 | 社会 |
| 6 政府 | 実現 | 対策 | 全力 | 日本 | 国民 |
| 7官 | 無駄 | 振興 | 選挙 | 共産党 | 保陣 |
| 8推進 | 日本 | 充実 | 国政 | 献金 | 安全 |
| 9 民 | 増税 | 促進 | 作り | 金権 | 地域 |
| 10 自民党 | AI JA | 安定 | 横浜 | 充 | 拉致 |
| 11 日本 | 一元化 | 確立 | 対策 | 選挙 | 経済 |
| 12 制度 | 政権 | 企業 | 中小 | 禁止 | 守る |
| 13 民間 | 子供 | 実現 | 発電 | 憲法 | Pay 20 |
| 14 年金 | 地域 | 中小 | 推進 | 腐敗 | 北東府豊 華 |
| 15 実現 | ひと | 育成 | エネルギー | 団体 | 教育 |
| 16 進める | サラリーマン | 制度 | 企業 | 区 | 責任 |
| 17 断行 | 制度 | 政治 | 声 | ソ連 | カ |
| 18 地方 | 級員 | 地域 | 実現 | 守る | 割る |
| 19 止める | 金 | 名篇 7止 | 活性 | 平和 | 安心 |
| 20 保障 | 民主党 | 事業 | 自民党 | 円 | 目指す |
| 21 財政 | 年間 | 20.英 | 地方 | 反対 | 89 9 |
| 22 作る | 一掃 | 確保 | 尽くす | Ti. | 憲法 |
| 23 贊成 | 郵政 | 強化 | 商店 | 是正 | 可能 |
| 24 社会 | 道路 | 教育 | いかす | 18 | in |
| 25 国民 | 交代 | 施設 | 全国 | 悪政 | 未来 |
| 26 公務員 | 社会保険庁 | 生活 | 政党 | 抜本 | 20 |
| 27 カ | 月額 | 支援 | 25 | 定数 | 再生 |
| 28 経済 | 手当 | 環境 | 支援 | 政党 | 将来 |
| 29 🖹 | 談合 | 発展 | 経済 | 金丸 | 解決 |
| 30 年心 | 专籍 | 练等 | 2至 2小 | 26 軍 | 其士 |

Change in proportion of 'Pork' Topic

Change in proportion of 'Pork' Topic



Change in proportion of 'Foreign Policy' Topic

Change in proportion of 'Foreign Policy' Topic

