



THE UNIVERSITY
of EDINBURGH

Text Technologies for Data Science

INFR11145

Comparing Text Corpora (2)

Instructor:
Björn Ross

06-Nov-2024

1

LDA Overview

Björn Ross, TTDS 2024/2025



THE UNIVERSITY
of EDINBURGH

2

Background: Plate Notation

3

Björn Ross, TTDS 2024/2025

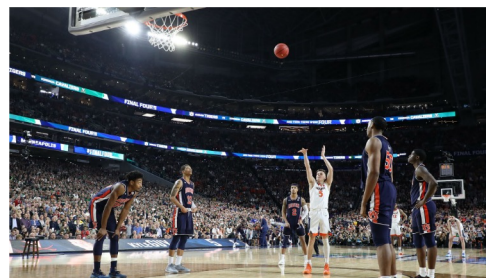


THE UNIVERSITY
of EDINBURGH

3

Background: Plate Notation

Make a
basket



4

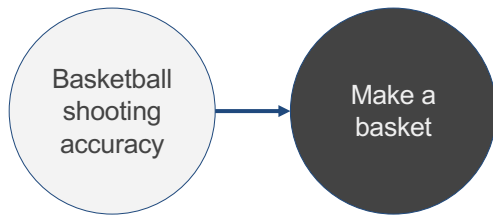
Björn Ross, TTDS 2024/2025



THE UNIVERSITY
of EDINBURGH

4

Background: Plate Notation

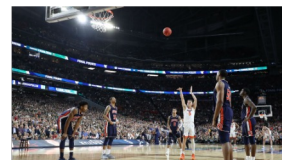
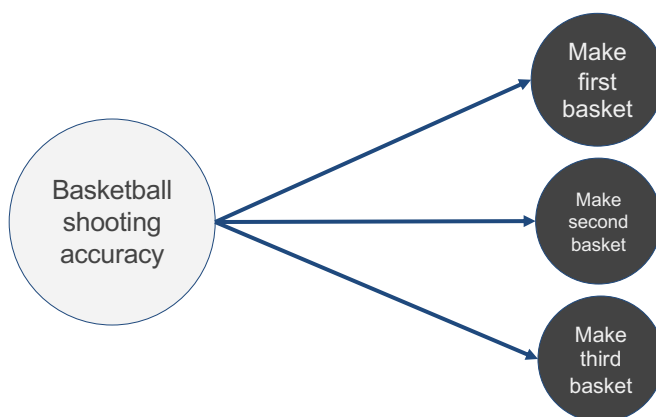


5

Björn Ross, TTDS 2024/2025

5

Background: Plate Notation

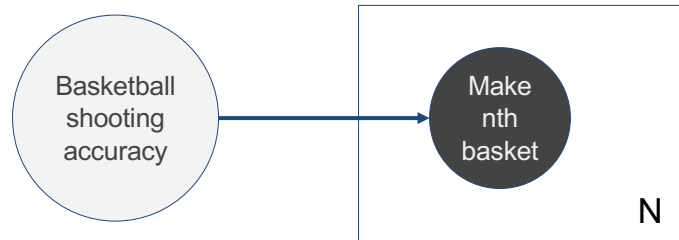


6

Björn Ross, TTDS 2024/2025

6

Background: Plate Notation



7

7

Latent Dirichlet Allocation

- Let's start with a very simple model
- We will work our way up to the full LDA model

8

Unigram Model

w is a word
 N words in a document

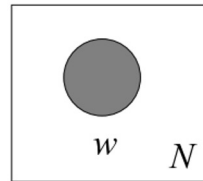


Figure from
Blei et al 2003

9

9

Unigram Model

w is a word
 N words in a document
 M documents in a corpus

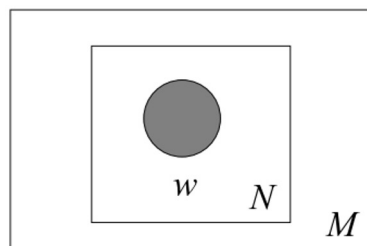


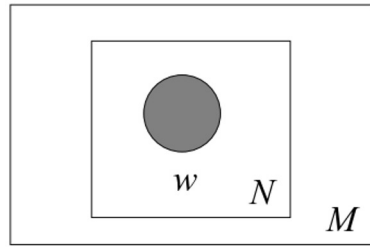
Figure from
Blei et al 2003

10

10

Unigram Model

w is a word
 N words in a document
 M documents in a corpus
 \mathbf{w} is a vector of words (i.e. doc)



$$p(\mathbf{w}) = \prod_{n=1}^N p(w_n)$$

Figure from
Blei et al 2003

11



11

Probability with a Unigram Model

$$p(\mathbf{w}) = \prod_{n=1}^N p(w_n)$$

What is the probability of the example sentence?

“My dog barked at another dog.”

word	my	at	dog	another	barked
probability	.10	.10	.05	.04	.03

12



12

Unigram Model...

- What is the point of making these models more complex?
- Why not just use the basic unigram model for everything?
- Remember:
 - Higher text probability *doesn't imply a better model*
 - We want to *accurately describe* the data
 - → higher probability for *real* documents, lower probability for noise

Mixture of Unigrams Model

z is the topic of a document

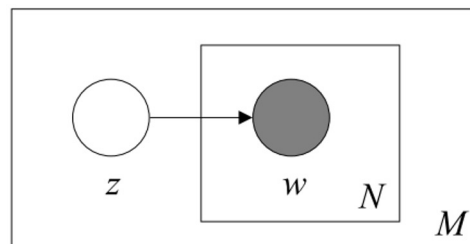
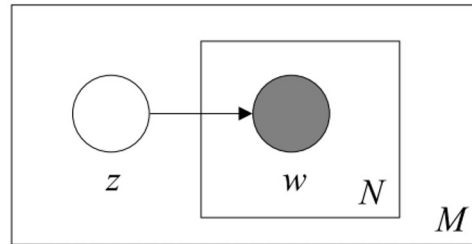


Figure from
Blei et al 2003

Mixture of Unigrams Model

z is the topic of a document



$$p(\mathbf{w}) = \sum_z p(z) \prod_{n=1}^N p(w_n | z)$$

Figure from
Blei et al 2003

16

16

Probability with Mixture of Unigrams

$$p(\mathbf{w}) = \sum_z p(z) \prod_{n=1}^N p(w_n | z).$$

What is the probability of the sentence?

Ignore stopwords: “my”, “after”, “the”

“My dog chased after the bus.”

w_i	cat	dog	chased	car	bus
$P(w_i z = \text{pets})$.20	.30	.10	.01	.01
$P(w_i z = \text{vehicles})$.01	.01	.10	.30	.20

$$p(z = \text{pets}) = 0.6,$$

$$p(z = \text{vehicles}) = 0.4$$

17

17

Probabilistic Latent Semantic Indexing

d is a document ID

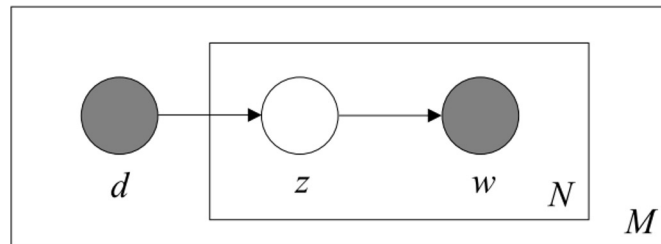


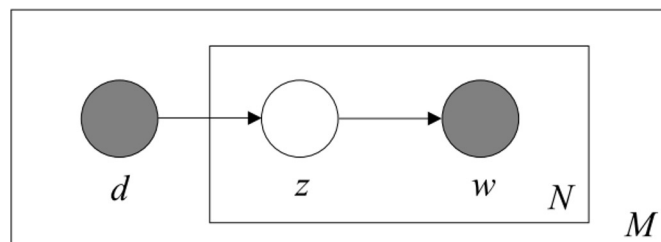
Figure from
Blei et al 2003

19

19

Probabilistic Latent Semantic Indexing

d is a document ID



$$p(d, w_n) = p(d) \sum_z p(w_n | z) p(z | d)$$

Figure from
Blei et al 2003

20

20

Probability with pLSI

w_i	cat	sat	down	car	broke
$p(w_i z = t_1)$.2	.1	.05	.01	.1
$p(w_i z = t_2)$.01	.05	.1	.3	.1

d_1 "The **cat** sat down."

$p(d = d_1)$.01
$p(z = t_1 d = d_1)$.6
$p(z = t_2 d = d_1)$.4

w_i	cat	sat	down	car	broke
$p(w_i z = t_1)$.2	.1	.05	.01	.1
$p(w_i z = t_2)$.01	.05	.1	.3	.1

What is the joint probability of the document and the word "cat"?²¹

Björn Ross, TTDS 2024/2025



21

Probability with pLSI

$$p(d, w_n) = p(d) \sum_z p(w_n | z) p(z | d)$$

Solution:

The **cat** sat down.

$$0.01 * (0.2 * 0.6 + 0.01 * 0.4) = 0.00124$$

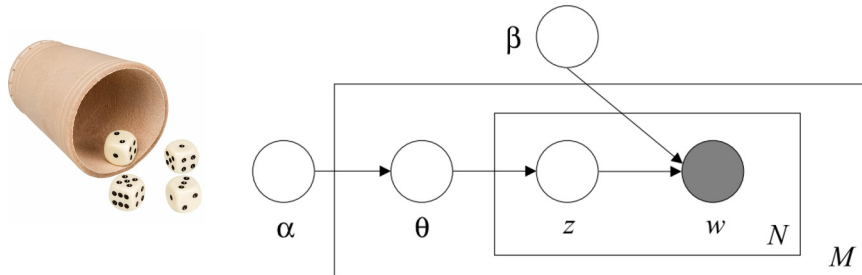
22

Björn Ross, TTDS 2024/2025



22

Latent Dirichlet Allocation

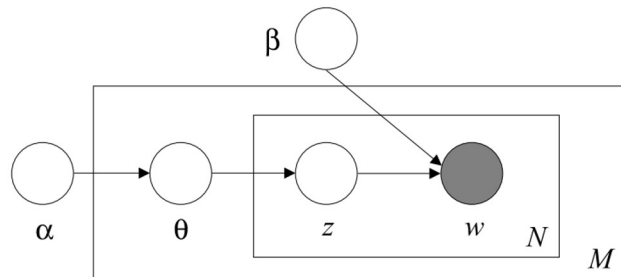


θ is the distribution over topics in a document
 α is the parameter of a Dirichlet distribution giving possible topic distributions within documents
 β gives word distributions within topics

Figure from Blei et al 2003

23

Latent Dirichlet Allocation



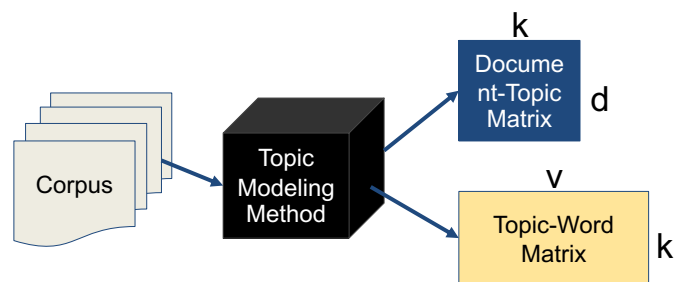
$$p(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

Figure from Blei et al 2003

24

Model Inference

- Want to learn the model parameters
- Exact inference becomes intractable



26

26

Model Inference

- Instead, use an approximate method such as:
 - Gibbs sampling
 - Variational Inference

27

27

Gibbs Sampling for LDA

Goal: Learn Φ , θ given a set of documents D

Φ = topic-word probabilities

θ = document-topic probabilities

Known:

corpus, α , β and the probability that a word is from a topic conditional on the assignments of all other words to topics

$$P(z_i = j | \mathbf{z}_{-i}, w_i, d_i, \cdot) \propto \frac{C_{w_i j}^{WT} + \beta}{\sum_{w=1}^W C_{w j}^{WT} + W\beta} \frac{C_{d_i j}^{DT} + \alpha}{\sum_{t=1}^T C_{d_i t}^{DT} + T\alpha}$$

Note: the \propto symbol means "proportional to"

28



28

Gibbs Sampling for LDA

Want to learn Φ , θ given a set of documents D

1. Assign each word a topic randomly
2. Calculate count matrices
3. Repeat until convergence:
 - For every document d
 - For every word i
 - Decrement count matrices C^{WT} and C^{DT} for current topic assignment
 - Sample a new topic assignment
 - Increment count matrices C^{WT} and C^{DT} for new topic assignment
4. Calculate Φ and θ

29



29

Gibbs Sampling for LDA

d1 Green eggs and ham.
 d2 Ham and green peppers.
 d3 Ham and cheese.

d1 Green eggs and ham.
 d2 Ham and green peppers.
 d3 Ham and cheese.

Random
 initialization.

30

30

Gibbs Sampling for LDA

C^{WT}	green	eggs	and	ham	peppers	cheese
t1	1	1	1	1	1	1
t2	1	0	2	2	0	0

Green eggs and ham.
 Ham and green peppers.
 Ham and cheese.

C^{DT}	d1	d2	d3
t1	2	2	2
t2	2	2	1

31

31

Gibbs Sampling for LDA

Assume (for the moment) $\alpha = \beta = 0$

θ	green	eggs	and	ham	peppers	cheese
t1	0.17	0.17	0.17	0.17	0.17	0.17
t2	0.20	0.00	0.40	0.40	0.00	0.00

Green eggs and ham.
Ham and green peppers.
Ham and cheese.

ϕ	d1	d2	d3
t1	0.50	0.50	0.66
t2	0.50	0.50	0.33

32

Björn Ross, TTDS 2024/2025



32

Gibbs Sampling for LDA

C^{WT}	green	eggs	and	ham	peppers	cheese
t1	1	1	1	1	1	1
t2	1	0	2	2	0	0

Green eggs and ham.
Ham and green peppers.
Ham and cheese.

C^{DT}	d1	d2	d3
t1	2	2	2
t2	2	2	1

33

Björn Ross, TTDS 2024/2025



33

Gibbs Sampling for LDA

C^{WT}	green	eggs	and	ham	peppers	cheese
t1	1	1	1	1	1	1
t2	1	0	2	2	0	0

Green eggs and ham.
Ham and green peppers.
Ham and cheese.

C^{DT}	d1	d2	d3
t1	2	2	2
t2	2	2	1

34

34

Gibbs Sampling for LDA

$$\frac{C_{w_{ij}}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \frac{C_{d_{ij}}^{DT} + \alpha}{\sum_{t=1}^T C_{d_{it}}^{DT} + T\alpha}$$

Assume (for the moment) $\alpha = \beta = 0$

C^{WT}	green	eggs	and	ham	peppers	cheese
t1	0	1	1	1	1	1
t2	1	0	2	2	0	0

Green eggs and ham.
Ham and green peppers.
Ham and cheese.

C^{DT}	d1	d2	d3
t1	1	2	2
t2	2	2	1

35

35

Gibbs Sampling for LDA

C^{WT}	green	eggs	and	ham	peppers	cheese
t1	0	1	1	1	1	1
t2	2	0	2	2	0	0

Green eggs and ham.
 Ham and green peppers.
 Ham and cheese.

C^{DT}	d1	d2	d3
t1	1	2	2
t2	3	2	1

36

36

Gibbs Sampling for LDA

C^{WT}	green	eggs	and	ham	peppers	cheese
t1	0	1	1	1	1	1
t2	2	0	2	2	0	0

Green eggs and ham.
 Ham and green peppers.
 Ham and cheese.

C^{DT}	d1	d2	d3
t1	1	2	2
t2	3	2	1

37

37

Gibbs Sampling for LDA

$$\frac{C_{w_{ij}}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \frac{C_{d_{ij}}^{DT} + \alpha}{\sum_{t=1}^T C_{d_{it}}^{DT} + T\alpha}$$

Assume (for the moment) $\alpha = \beta = 0$

C^{WT}	green	eggs	and	ham	peppers	cheese
t1	0	0	1	1	1	1
t2	2	0	2	2	0	0

Green eggs and ham.
 Ham and green peppers.
 Ham and cheese.

C^{DT}	d1	d2	d3
t1	0	2	2
t2	3	2	1

38

Björn Ross, TTDS 2024/2025

38

Gibbs Sampling for LDA

$$\frac{C_{w_{ij}}^{WT} + \beta}{\sum_{w=1}^W C_{wj}^{WT} + W\beta} \frac{C_{d_{ij}}^{DT} + \alpha}{\sum_{t=1}^T C_{d_{it}}^{DT} + T\alpha}$$

$C^{WT} + \alpha$	green	eggs	and	ham	peppers	cheese
t1	0.01	0.01	1.01	1.01	1.01	1.01
t2	2.01	0.01	2.01	2.01	0.01	0.01

Green eggs and ham.
 Ham and green peppers.
 Ham and cheese.

$C^{DT} + \beta$	d1	d2	d3
t1	0.01	2.01	2.01
t2	3.01	2.01	1.01

39

Björn Ross, TTDS 2024/2025

39

Gibbs Sampling for LDA

- Repeat until convergence
- Probabilistic algorithm – results depend on random initialisation and random samples!



Topic Modeling Examples



What do students look for in a professor?

Topic	Sample words
Approachability	prof, fair, clear, helpful, teaching, approachable, nice, organized, extremely, friendly, super, amazing
Clarity	understand, hard, homework, office, material, clear, helpful, problems, explains, accent, questions, extremely
Course Logistics	book, study, boring, extra, nice, credit, lot, hard, attendance, make, fine, attention, pay, mandatory
Enthusiasm	teaching, passionate, awesome, enthusiastic, professors, loves, cares, wonderful, fantastic, passion
Expectations	hard, work, time, lot, comments, tough, expects, worst, stuff, avoid, horrible, classes
Helpfulness	helpful, nice, recommend, cares, super, understanding, kind, extremely, effort, sweet, friendly, approachable
Humor	guy, funny, fun, awesome, cool, entertaining, humor, hilarious, jokes, stories, love, hot, enjoyable
Interestingness	interesting, material, recommend, lecturer, engaging, classes, knowledgeable, enjoyed, loved, topics
Readings/ Discussions	readings, papers, writing, ta, interesting, discussions, grader, essays, boring, books, participation
Study Material	exams, notes, questions, material, textbook, hard, slides, study, answer, clear, tricky, attend, long, understand

Azab, Mihalcea, and Abernathy, 2016

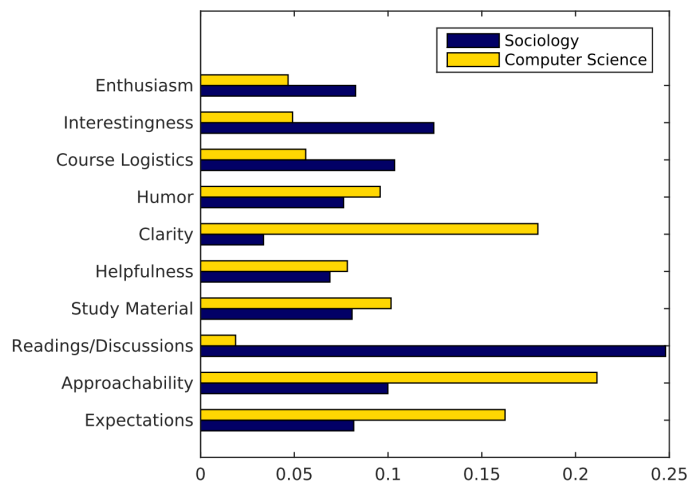
Björn Ross, TTDS 2024/2025



THE UNIVERSITY
of EDINBURGH

42

What do students look for in a professor?



Azab, Mihalcea, and Abernathy, 2016

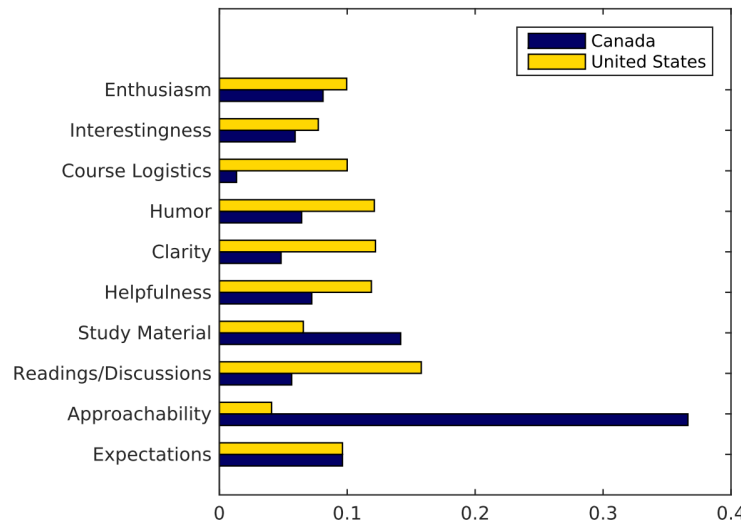
Björn Ross, TTDS 2024/2025



THE UNIVERSITY
of EDINBURGH

43

What do students look for in a professor?



Azab, Mihalcea, and Abernathy, 2016

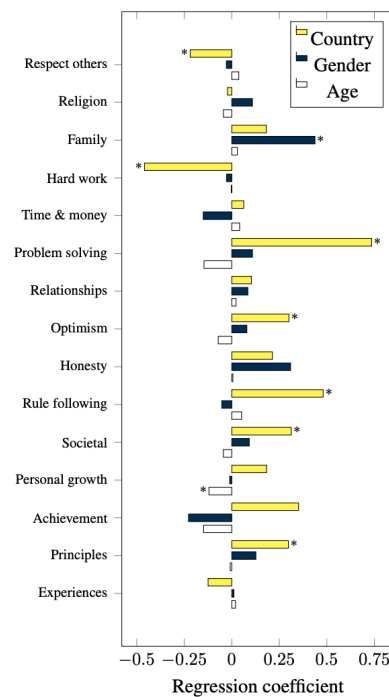
Björn Ross, TTDS 2024/2025



44

How do personal attributes relate to values?

Theme	Example Words
Respect others	people, respect, care, human, treat
Religion	god, heart, belief, religion, right
Family	family, parent, child, husband, mother
Hard Work	hard, work, better, honest, best
Time & Money	money, work, time, day, year
Problem solving	consider, decision, situation, problem
Relationships	family, friend, relationship, love
Optimism	enjoy, happy, positive, future, grow
Honesty	honest, truth, lie, trust, true
Rule following	moral, rule, principle, follow
Societal	society, person, feel, thought, quality
Personal Growth	personal, grow, best, decision, mind
Achievement	heart, achieve, complete, goal
Principles	important, guide, principle, central
Experiences	look, see, experience, choose, feel



Wilson, Mihalcea, Boyd, and Pennebaker 2016

Björn Ross, TTDS 2024/2025

45

Annotation + Classification

Björn Ross, TTDS 2024/2025



46

Annotation + Classification

- Method 1: Traditional Supervised Learning
 - Annotate representative samples
 - Train a classifier
 - Apply to rest of data
- Method 2: Transfer Learning
 - Find another large, but similar dataset
 - Train a classifier on that dataset
 - *Optionally: fine-tune classifier to your smaller dataset*
 - Apply to rest of your data

Björn Ross, TTDS 2024/2025



47

After Classification

- Which features are most relevant for each class?
- What are common words/topics for each class?
- How do predicted classes relate to other variables?
- *More about text classification coming up next week!*

Wrap-up

- Content analysis background
- Word-level differences
- Dictionaries and Lexica
- Topic modeling
- Annotation + classification

Readings

- [Manning: IR book](#) section 13.5
- [“Probabilistic Topic Models”](#) by David Blei
- [“Latent Dirichlet Allocation”](#) by David Blei, Andrew Y. Ng, and Michael I. Jordan
- [“Probabilistic Topic Models”](#) by Mark Steyvers and Tom Griffiths

To watch:

- [Guest lecture](#) (2017) by David Blei at University of Edinburgh School of Informatics