

Latent Dirichlet Allocation

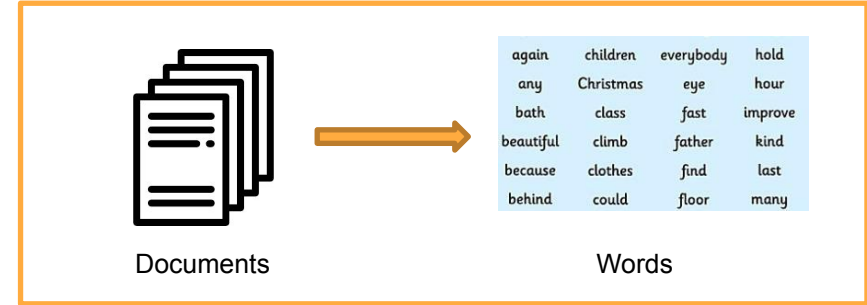
Latent Dirichlet Allocation

- Generative probabilistic model

Latent Dirichlet Allocation

- Generative probabilistic model

Finds topics from a corpus
Annotates documents with topics



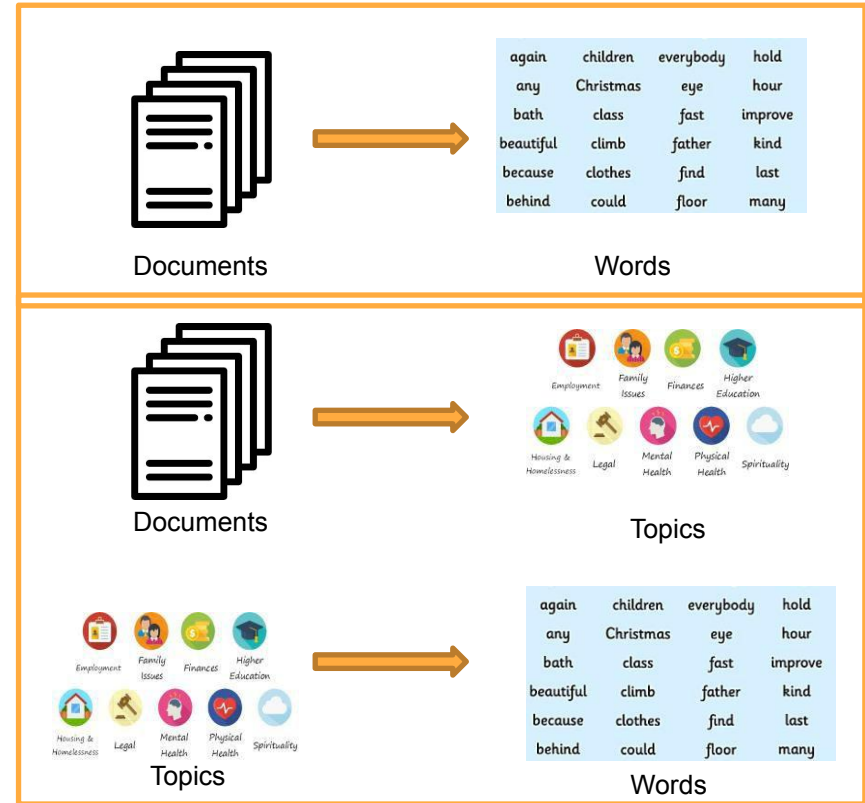
Latent Dirichlet Allocation

- Generative probabilistic model

Finds topics from a corpus
Annotates documents with topics

- LDA Assumptions

Documents = mixture of topics
Topics = mixture of words



Latent Dirichlet Allocation

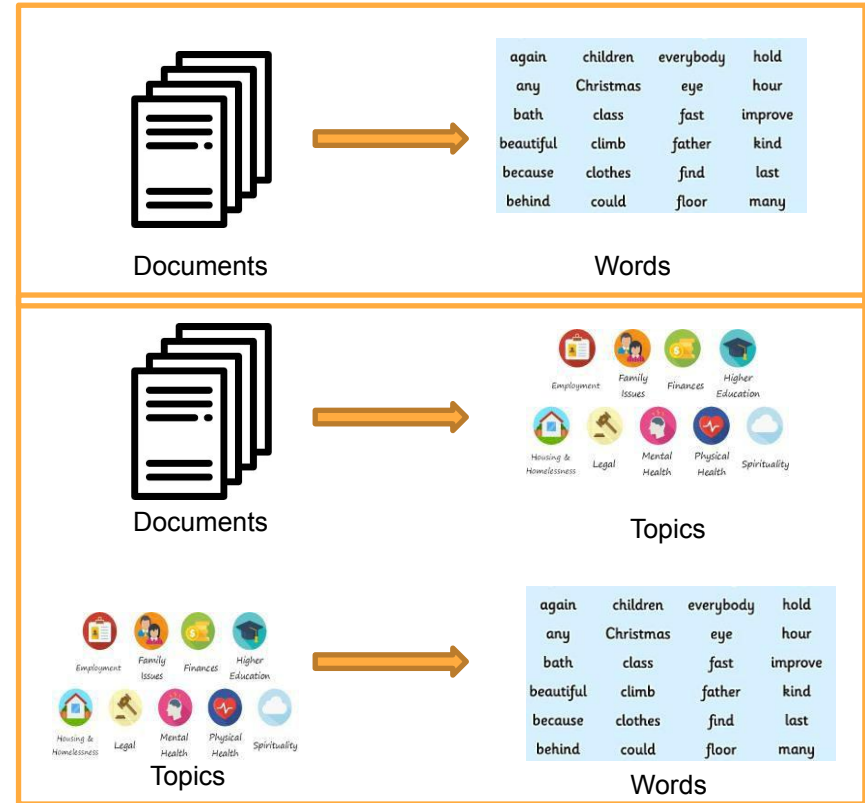
- Generative probabilistic model

Finds topics from a corpus
Annotates documents with topics

- LDA Assumptions

Documents = mixture of topics
Topics = mixture of words

- Documents : Probability Distributions of Topics
Topics : Probability Distributions of Words



Latent Dirichlet Allocation

Latent Dirichlet Allocation

- Corpus : Document Word Matrix
- Document Word Matrix = Document Topic Matrix + Topic Word Matrix

Latent Dirichlet Allocation

- Corpus : Document Word Matrix

Latent Dirichlet Allocation

- Corpus : Document Word Matrix
- Document Word Matrix = Document Topic Matrix + Topic Word Matrix

	W1	W2	W3	<u>W_n</u>
D1	0	2	1	3
D2	1	4	0	0
D3	0	2	3	1
<u>D_n</u>	1	1	3	0



	K1	K2	K3	K
D1	1	0	0	1
D2	1	1	0	0
D3	1	0	0	1
<u>D_n</u>	1	0	1	0

	W1	W2	W3	<u>W_m</u>
K1	0	1	1	1
K2	1	1	1	0
K3	1	0	0	1
K	1	1	0	0

Latent Dirichlet Allocation

- Corpus : Document Word Matrix
- Document Word Matrix = Document Topic Matrix + Topic Word Matrix

	W1	W2	W3	<u>W_n</u>
D1	0	2	1	3
D2	1	4	0	0
D3	0	2	3	1
<u>D_n</u>	1	1	3	0



	K1	K2	K3	K
D1	1	0	0	1
D2	1	1	0	0
D3	1	0	0	1
<u>D_n</u>	1	0	1	0

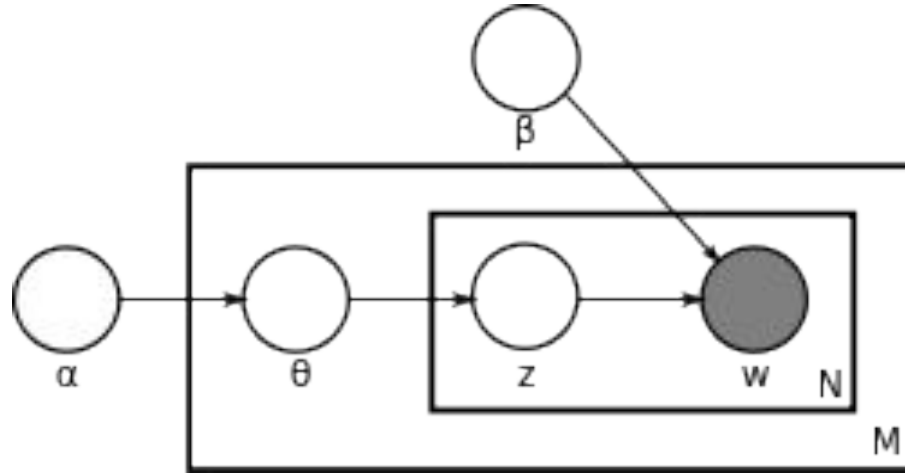
	W1	W2	W3	<u>W_m</u>
K1	0	1	1	1
K2	1	1	1	0
K3	1	0	0	1
K	1	1	0	0

- Goal – Optimize representations

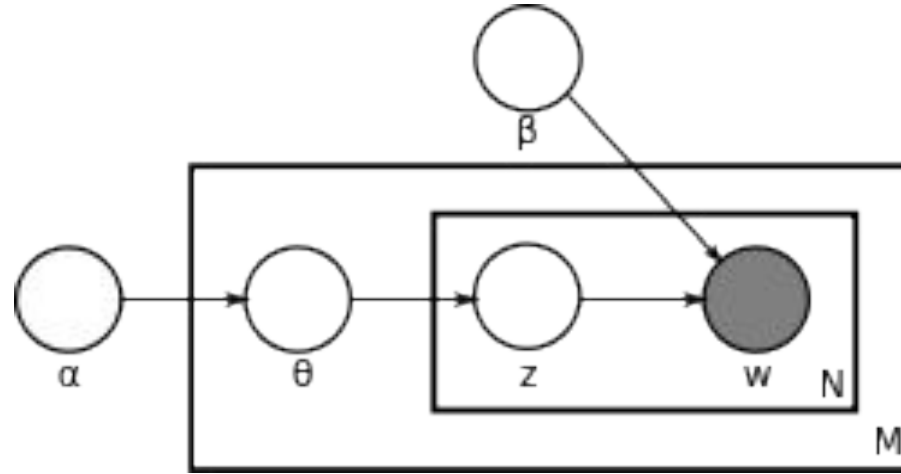
Document Topic distributions
Topic Terms distributions

Latent Dirichlet Allocation

Latent Dirichlet Allocation

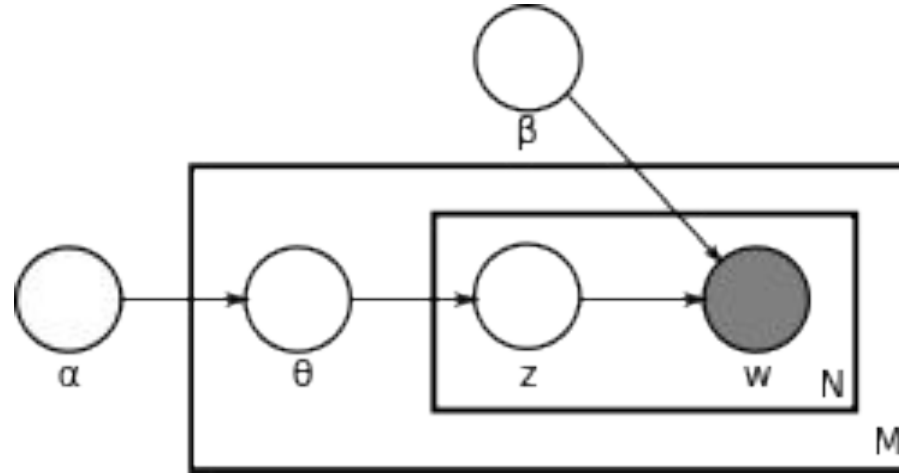


Latent Dirichlet Allocation



- M : Total Documents in Corpus
- N : No of words in a Document
- w : Word in a document
- z : Latent topic assigned to the word
- theta : Topic Distribution

Latent Dirichlet Allocation



- M : Total Documents in Corpus
 N : No of words in a Document
 w : Word in a document
 z : Latent topic assigned to the word
 θ : Topic Distribution
- Alpha, Beta – LDA model parameters

Latent Dirichlet Allocation

Latent Dirichlet Allocation

- Corpus:

$D1 = (w_1, w_2, w_3, w_4, \dots, w_n)$

$D2 = (w'_1, w'_2, w'_3, w'_4, \dots, w'_n)$

$D3 = (w''_1, w''_2, w''_3, w''_4, \dots, w''_n)$

...

...

$D_m = (w_1, w_2, w_3, w_4, \dots, w_n)$

Latent Dirichlet Allocation

- Corpus:

$D1 = (w_1, w_2, w_3, w_4, \dots, w_n)$

$D2 = (w'_1, w'_2, w'_3, w'_4, \dots, w'_n)$

$D3 = (w''_1, w''_2, w''_3, w''_4, \dots, w''_n)$

...

...

$D_m = (w_1, w_2, w_3, w_4, \dots, w_n)$

- First step : Assign random topics to each word

Latent Dirichlet Allocation

- Corpus:

D1 = (w1, w2, w3, w4, wn)

D2 = (w'1, w'2, w'3, w'4,
w'n)

D3 = (w''1, w''2, w''3, w''4, w''n)

...

...

Dm = (w1, w2, w3, w4, wn)

- First step : Assign random topics to each word

D1 = (w1 (k4), w2 (k2), w3 (k2), w4 (k2), wn (k3))

D2 = (w'1 (k1), w'2 (k7), w'3 (k3), w'4 (k6), w'n (k2))

D3 = (w''1(k5), w''2 (k4), w''3 (k1), w''4 (k5), w''n
(k1))

...

...

Dm = (w1 (k4), w2 (k2), w3 (k6), w4 (k1), wn (k2))

Latent Dirichlet Allocation

$D1 = (w_1 (k_4), w_2 (k_2), w_3 (k_2), w_4 (k_2), \dots, w_n (k_3))$
 $D2 = (w'_1 (k_1), w'_2 (k_7), w'_3 (k_3), w'_4 (k_6), \dots, w'_n (k_2))$
 $D3 = (w''_1 (k_5), w''_2 (k_4), w''_3 (k_1), w''_4 (k_5), \dots, w''_n (k_1))$

Documents : Mixture of Topics:

$D1 = k_4 + k_2 + k_2 + k_2 + \dots$
 $k_3 \quad D2 = k_1 + k_7 + k_3 + k_6 +$
 $\dots k_2 \quad D3 = k_5 + k_4 + k_1 + k_5$
 $+ \dots k_1 \quad Dn = \dots$
 $.$

Latent Dirichlet Allocation

$D1 = (w_1 (k_4), w_2 (k_2), w_3 (k_2), w_4 (k_2), \dots, w_n (k_3))$
 $D2 = (w'_1 (k_1), w'_2 (k_7), w'_3 (k_3), w'_4 (k_6), \dots, w'_n (k_2))$
 $D3 = (w''_1 (k_5), w''_2 (k_4), w''_3 (k_1), w''_4 (k_5), \dots, w''_n (k_1))$

Documents : Mixture of Topics:

$D1 = k_4 + k_2 + k_2 + k_2 + \dots$
 $k_3 \quad D2 = k_1 + k_7 + k_3 + k_6 + \dots$
 $k_2 \quad D3 = k_5 + k_4 + k_1 + k_5 + \dots$
 $k_1 \quad Dn = \dots$

Topics : Mixture of Terms:

$k_1 = w'_1 + w''_3$

$k_2 = w_2 + w_3 + w_4 + \dots$

\dots

$k_n = w_i + \dots$

Latent Dirichlet Allocation

Optimization Steps:

Iterate : each document d

Iterate : each word w

Latent Dirichlet Allocation

Optimization Steps:

Iterate : each document d

Iterate : each word w

-Assume that all topic assignments except the current word are correct

Latent Dirichlet Allocation

Optimization Steps:

Iterate : each document d

Iterate : each word w

-Assume that all topic assignments except the current word are correct

-compute p_1, p_2

Latent Dirichlet Allocation

Optimization Steps:

Iterate : each document d

Iterate : each word w

-Assume that all topic assignments except the current word are correct

-compute $p1$, $p2$

$p1 = \text{proportion (topic } t \text{ / document } d)$ $p2 = \text{proportion (word } w \text{ / topic } t)$

$p1$ -> proportion of **words in document d that are currently assigned to topic t**

$p2$ -> proportion of **assignments to topic t that come from w , over all documents**

Latent Dirichlet Allocation

- Reassign word w of document d a new topic k'
 - Where we choose topic k' with a new probability = $p_1 * p_2$

Latent Dirichlet Allocation

- Reassign word w of document d a new topic k'
 - Where we choose topic k' with a new probability = $p_1 * p_2$
- Repeated large number of times until steady state

Thank You