

# Text Mining: Latent Dirichlet Allocation

Universidad de los Andes

Juan Sebastián Moreno Pabón

3 de septiembre de 2019

# Modelos de tópicos

- El objetivo de un modelo de tópicos es descubrir las temáticas subyacentes en un corpus de texto.
- El objetivo es muy similar al análisis de clústeres pero para textos.
- Retos:
  - Un texto puede referirse a múltiples temáticas.
  - ¿Cómo tratar un texto como data estructurada?

# Latent Dirichlet Allocation (LDA)

<b>Title</b> word	<b>Title</b> word	<b>Title</b> word
<b>Title</b> word	<b>Title</b> word	<b>Title</b> word
<b>Title</b> word	<b>Title</b> word	<b>Title</b> word
<b>Title</b> word	<b>Title</b> word	<b>Title</b> word

remove if word appears  
in more than 10% of  
the articles

keep top n words  
Recommended 50,000 - 100,000

discard the rest

remove if the word  
appears in less  
than 20 articles

## Let the right words in

### Word length

- 1: I
- 2: do, be, am, ...
- 3: ice, was, who, ...

- 16: videoconferences, ...
- 17: superbillionaires, ...
- 18: intellectualization, ...

### Stoptlists:

general terms  
last names, first names  
countries, cities

### Lemmatization

am, are, is = be

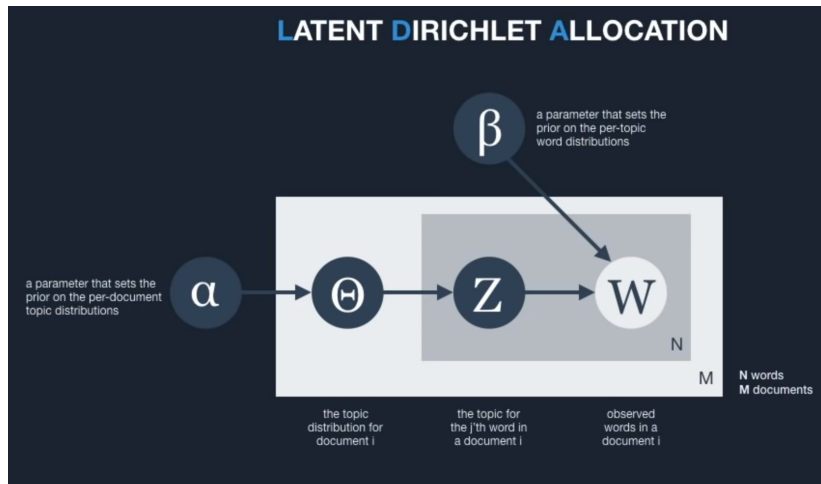
### Parts of Speech:

**NN** - noun (computer, car, cake, ...)  
**VB** - verb (play, install, commit, ...)  
**RB** - adverb (today, quickly, patiently, ...)  
**JJ** - adjective (red, awesome, big, ...)  
**IN** - preposition (of, about, from, ...)

# Latent Dirichlet Allocation (LDA)

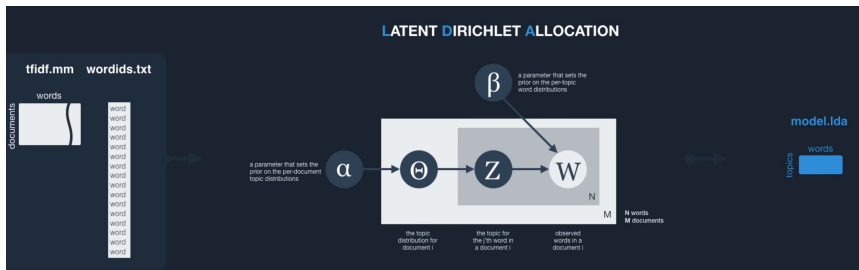


# Latent Dirichlet Allocation (LDA)

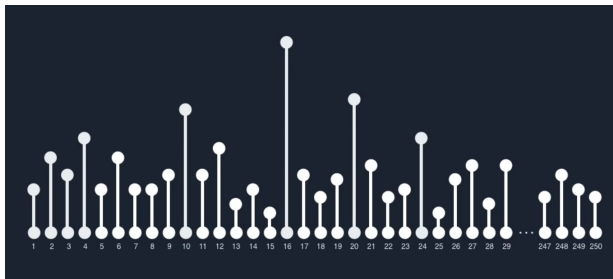


¿Qué hace Alfa y Beta?

# LDA



# LDA





# LDA

