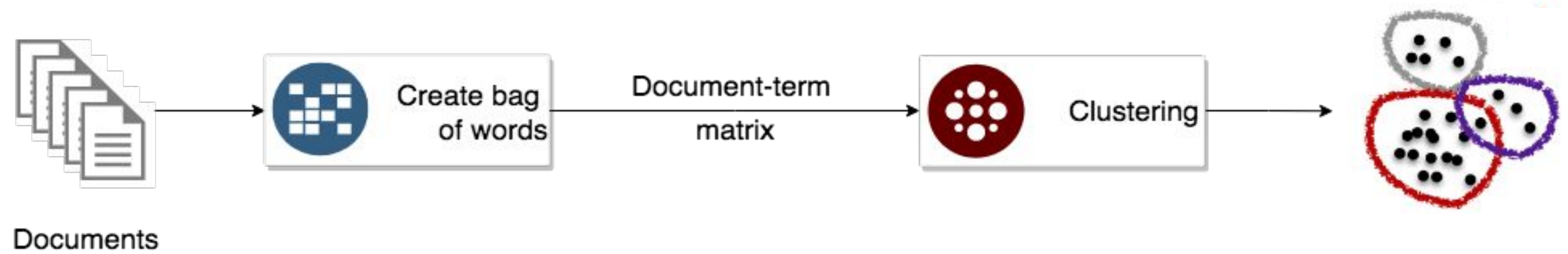
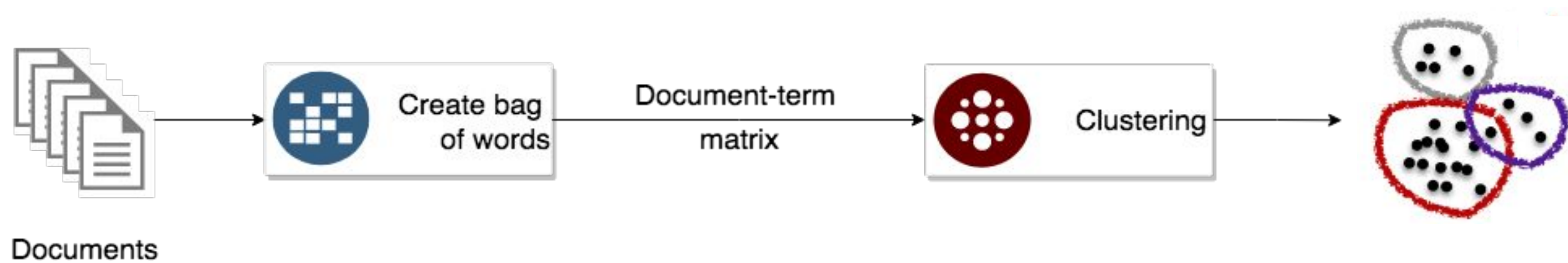


Análise de Mídias Sociais e Mineração de Texto

Modelagem de Tópicos

Laura de Oliveira F. Moraes





documents

	<i>for</i>	<i>if</i>	<i>append</i>
<i>1</i>	<i>2</i>	<i>1</i>	<i>1</i>
<i>2</i>	<i>0</i>	<i>1</i>	<i>1</i>
<i>3</i>	<i>1</i>	<i>2</i>	<i>0</i>

terms/tokens/words



Separa o texto em:

N = 1 : This is a sentence *unigrams:*

N = 2 : This is a sentence *bigrams:*

N = 3 : This is a sentence *trigrams:*

- Entre 5% a 50%
- Remove palavras comuns como preposições, artigos, etc

Como contar:

1. Uma vez por documento **OU**
2. Cada ocorrência no documento

Normalização da contagem:

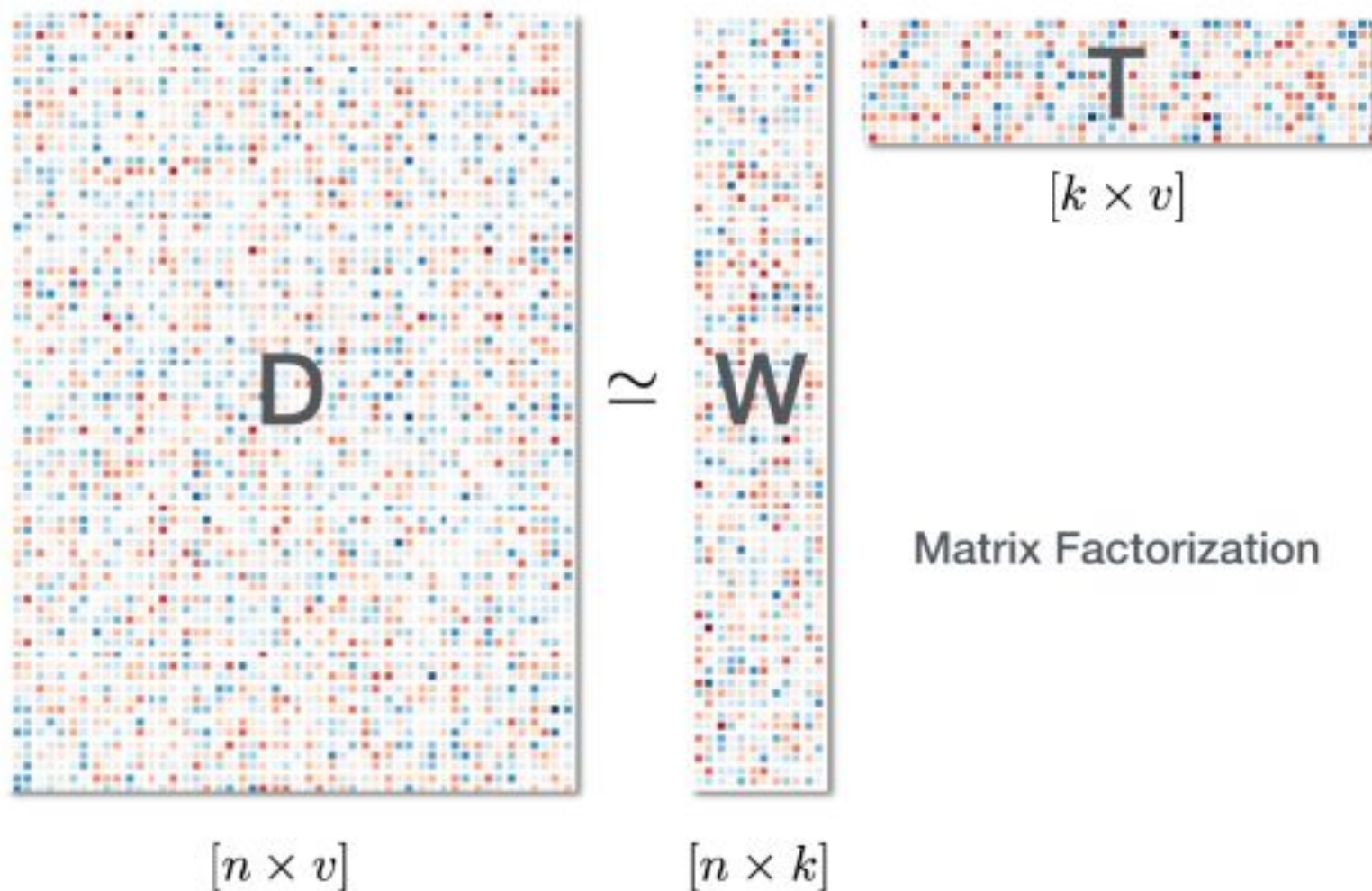
1. Contagem Normal **OU**
2. TF-IDF



Methods for Topic Modeling

1. Non-negative Matrix Factorization (NMF)

(Lee et al 1999, Cichocki et al 2009)



When constructing the document-term matrix, the terms counts can be interpreted as a set of the visible variables generated from an underlying set of hidden variables (topics). With this in mind, we can model the (hidden) topics as linear combinations of terms (a sum of each term weighted by its corresponding importance) and the documents as a weighted-sum of the topics it belongs to.

This factorization has a special property of only allowing non-negative values in its entries. By imposing this constraint, the resulting factorized matrices can be described essentially as weighted-union of sets, which is well-suited for human interpretability.



Clustering

Methods for Topic Modeling

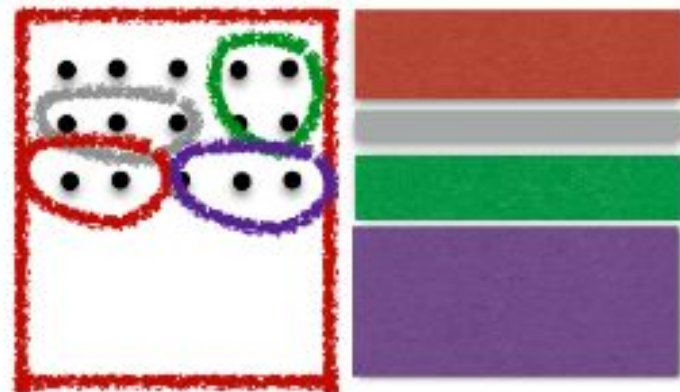
2. Latent Dirichlet Allocation (LDA)

(Blei et al 2003. Steyvers et al 2010)

LDA is a generative probabilistic model that describes how documents in a collection are created. It assumes that:

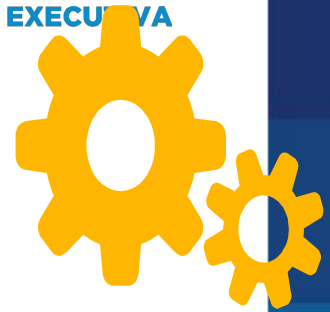
1. A document is a mixture of topics

Document



• words/tokens

distribution of topics



Clustering

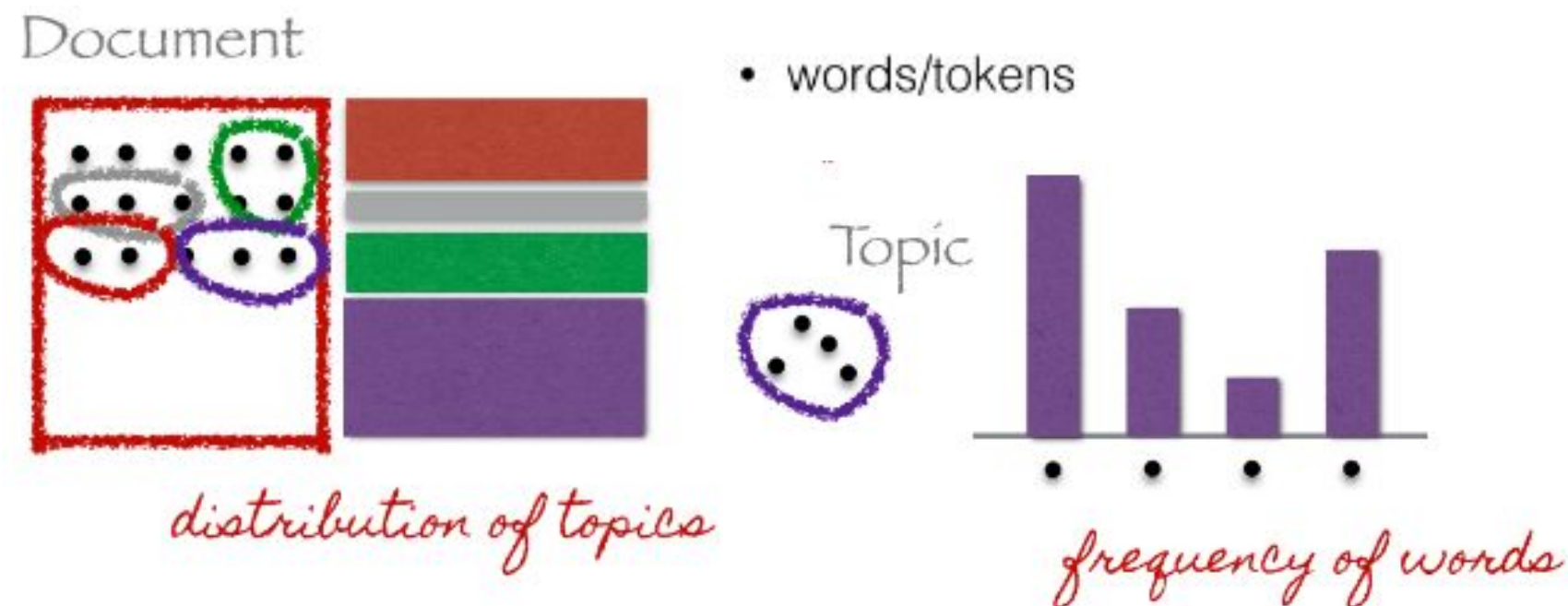
Methods for Topic Modeling

2. Latent Dirichlet Allocation (LDA)

(Blei et al 2003. Steyvers et al 2010)

LDA is a generative probabilistic model that describes how documents in dataset were created. It assumes that:

1. A document is a mixture of topics
2. A topic is a distribution over words

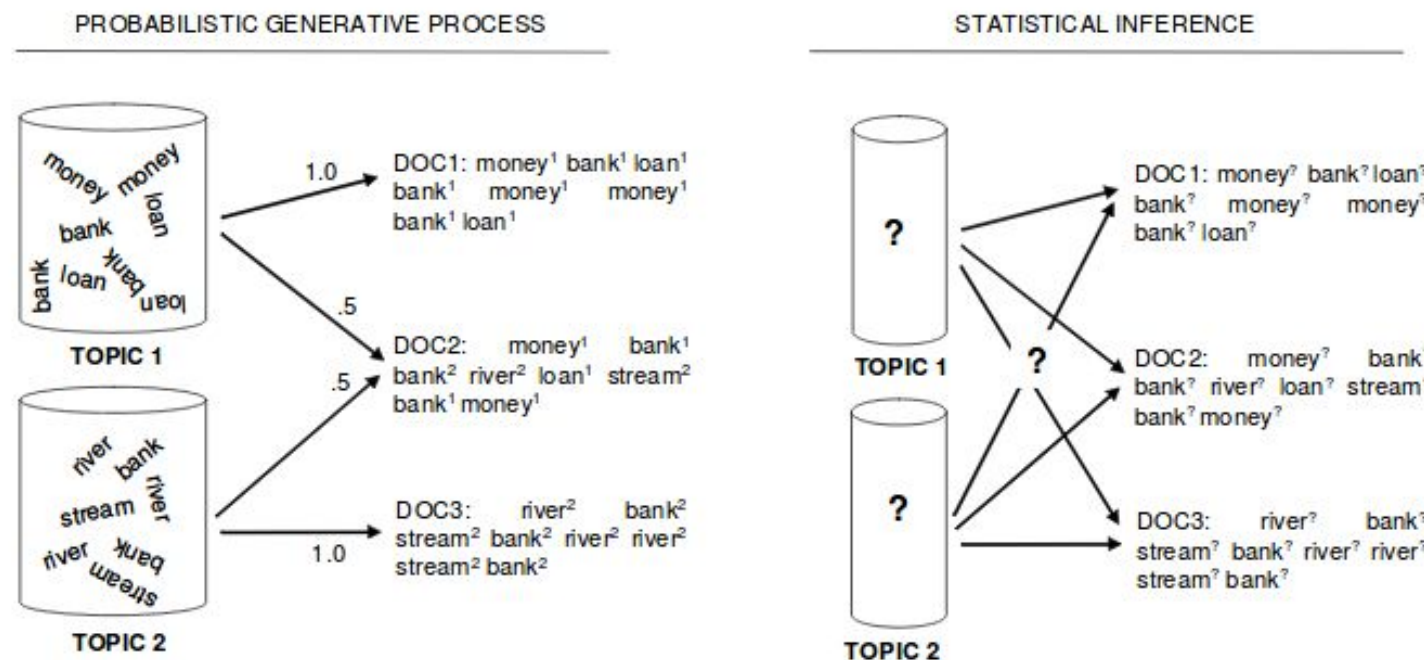




Methods for Topic Modeling

2. Latent Dirichlet Allocation (LDA)

(Blei et al 2003. Steyvers et al 2010)



To generate a document, you simply sample from the distributions:

1. Sample a topic for the document
2. Sample a word from the topic

Doing this iteratively, you will generate a document.

However, once you have a dataset, meaning a collection of documents, what we need is to discover these distributions. The LDA algorithm tries to backtrack this probabilistic model to find a set of topics that are likely to have generated the collection.

TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

Term frequency

Inverse document frequency

Number of times term t appears in a doc, d

$$\log \frac{1 + n}{1 + \text{df}(d, t) + 1}$$

of documents

Document frequency of the term t