

## Topic Models

Tecnologías de búsqueda en la web

Marcelo Mendoza



## Topic Models: LDA

MM

INF-335

1 / 7

MM

INF-335

2 / 7

Topic Models

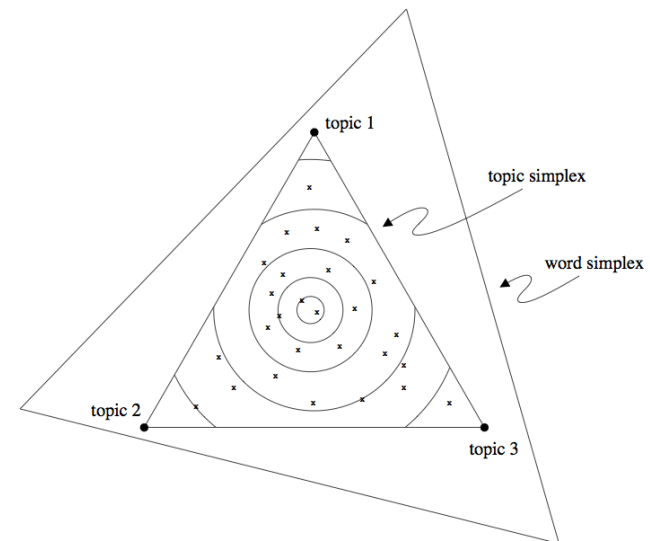
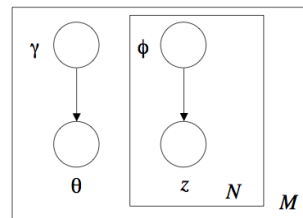
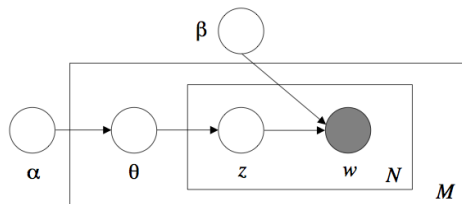
Topic Models

Topic Models

Topic Models

## Decoupling trick

## LDA: Interpretación geométrica



MM

INF-335

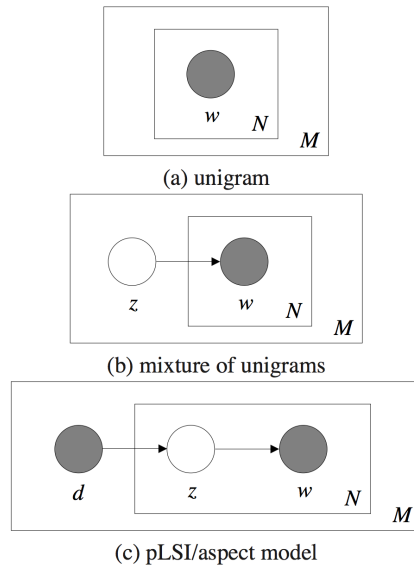
3 / 7

MM

INF-335

4 / 7

## Modelos Gráficos



MM

INF-335

5 / 7

MM

INF-335

6 / 7

## Gensim (<http://radimrehurek.com/gensim/>)

### Installing:

```
1 > sudo pip install --upgrade gensim
```

### Corpus (preprocesamiento):

```
1 > python
2 >>> from gensim import corpora, models
3 >>> documents = ["Human machine interface for lab abc computer applications",
4 ...             "A survey of user opinion of computer system response time",
5 ...             "The EPS user interface management system",
6 ...             "System and human system engineering testing of EPS",
7 ...             "Relation of user perceived response time to error measurement",
8 ...             "The generation of random binary unordered trees",
9 ...             "The intersection graph of paths in trees",
10 ...            "Graph minors IV Widths of trees and well quasi ordering",
11 ...            "Graph minors A survey"]
12 >>> import nltk
13 >>> from nltk.corpus import stopwords
14 >>> stop = stopwords.words('english')
15 >>> texts = [[word for word in document.lower().split() if word not in stop]
16 ...          for document in documents]
```

## Gensim

### Corpus (representación BOW):

```
1 >>> dictionary = corpora.Dictionary(texts)
2 >>> corpus = [dictionary.doc2bow(text) for text in texts]
3 >>> corpora.MmCorpus.serialize('/tmp/deerwester.mm', corpus)
```

### LDA fitting:

```
1 >>> lda = models.LdaModel(corpus, id2word=dictionary, num_topics=2)
```

### Term dist. (print\_topics(num\_topics)):

```
1 >>> lda.print_topics(2)
```

### Top-terms per topic (print\_topic(topicid,topn)):

```
1 >>> lda.print_topic(0, topn=5)
2 >>> lda.print_topic(1, topn=5)
```

### Document dist.:

```
1 >>> doc2bow = dictionary.doc2bow(texts[1])
2 >>> lda[doc2bow]
```

MM

INF-335

7 / 7

### LDA update (sobre un nuevo corpus):

```
1 >> lda.update(new_corpus)
```