

Statistical Analysis of Text - a mini-course

Word embeddings and Topic models

Mattias Villani

Statistiska institutionen
Stockholms universitet

Institutionen för datavetenskap
Linköpings universitet



Overview Lecture 3

- Distributional semantics and **Word embeddings**
- **Topic models**
- Demo of **topicmodels package** in R

The distributional semantics hypothesis

- **Semantics** - the meaning of words
- **Distributional semantics** - categorizing meaning by the distributional properties in **large samples**
- The distributional semantics hypothesis
"a word is characterized by the company it keeps"
Firth (1957)

Distributional semantics

- Word meaning comes from **textual context**

"cold"

"It's cold outside."

"I'm having a cold"

"I'm cold"

- **Different contexts** (sentence, word windows, documents)
- **Different context size** - different properties
 - ▶ Short distance context, syntagmatic similarities
 - ▶ Long distance context, topical similarities

Co-occurrence matrix

A friend in need is a friend indeed.
She is my friend indeed.

	Doc 1	Doc 2
a	2	0
friend	2	1
in	1	0
indeed	1	1
is	1	1
my	0	1
need	1	0
she	0	1

Co-occurrence matrix

A friend in need is a friend indeed.

She is my friend indeed.

■ Context window of one step

	a	friend	in	indeed	is	my	need	she
a	2	2	0	0	1	0	0	0
friend	2	3	1	2	0	1	0	0
in	0	1	1	0	0	0	1	0
indeed	0	2	0	2	0	0	0	0
is	1	0	0	0	2	1	1	1
my	0	1	0	0	1	1	0	0
need	0	0	1	0	1	0	1	0
she	0	0	0	0	1	0	0	1

Word embedding

- Reduce co-occurrence matrix to a **lower dimension**
- Often part of more complex models
- Popular approaches (word-word)
 - ▶ **Random indexing**
 - ▶ **Word2Vec**
- Popular approaches (word-doc)
 - ▶ **Latent Semantic analysis** (SVD decomposition)
 - ▶ **Topic models**

Topic models

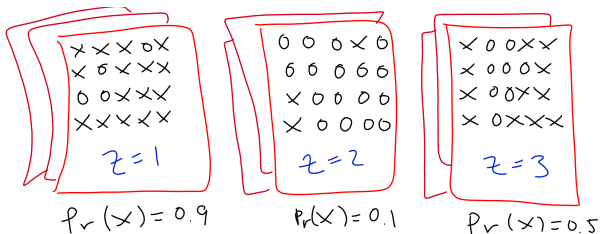
- Models for **unsupervised learning**, but more recently also for **supervised learning**.
- **Probabilistic generative** models.
- Very popular model in applications and research. > 12000 Google scholar citations in 11 years.
- **Extensions**: nGrams, supervised, nonparametric, relational topics, correlated topics, dynamically time-varying topics etc.
- The basic topic models are extensions of the bag-of-words (unigram) model.
- **Unigram model**: each word is assumed to be drawn from the same word (term) distribution.

$$\hat{P}(w) = \frac{\#w}{N}$$

Mixture of unigrams

■ Mixture of unigrams:

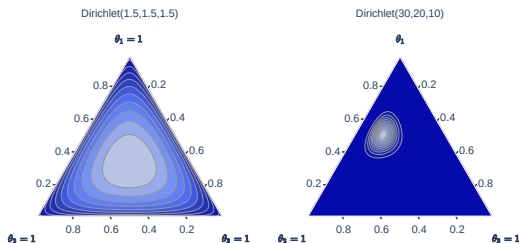
- 1 Draw a *topic* z_d for the d th document from a topic distribution $\theta = (\theta_1, \dots, \theta_K)$.
- 2 Conditional on the drawn topic z_d draw words from a word distribution for that topic.



- Topic models are **mixed-membership models**: each document can belong to **several topics simultaneously**.

Multinomial and Dirichlet distributions

- **Multinomial distribution**: discrete unordered $X \in \{1, 2, \dots, K\}$.
 - ▶ $Pr(X = k) = \theta_k$
 - ▶ Parameters $\theta = (\theta_1, \dots, \theta_K)$.
- **Dirichlet distribution**: random vector $X = (X_1, \dots, X_K)$.
 - ▶ **Unit simplex** $X_1 + X_2 + \dots + X_K = 1$
 - ▶ Parameters: $\alpha = (\alpha_1, \dots, \alpha_K)$
 - ▶ Uniform distribution: $\alpha = (1, 1, \dots, 1)$
 - ▶ Small variance (informative) when the α 's are large.
 - ▶ “Bathtub shape” when $\alpha_k < 1$ for all k .



Generating a corpus from a topic model

■ Assume that we have:

- ▶ A fixed vocabulary V
- ▶ D documents
- ▶ N words in each document
- ▶ K topics

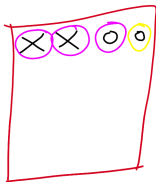
1 For each topic ($k = 1, \dots, K$):

- a. Draw a distribution over the words $\beta_k \sim \text{Dir}(\eta, \eta, \dots, \eta)$

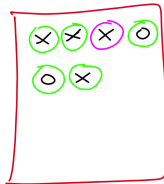
2 For each document ($d = 1, \dots, D$):

- a. Draw a vector of topic proportions $\theta_d \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$
- b. **For each word** ($n = 1, \dots, N$):
 - i. Draw a topic assignment $z_{d,n} \sim \text{Multinomial}(\theta_d)$
 - ii. Draw a word $w_{d,n} \sim \text{Multinomial}(\beta_{z_{d,n}})$

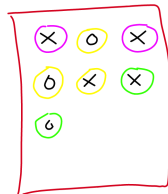
(horrible) picture of a topic model



$$\theta_1 = (\underline{0.9} \quad \underline{0.1} \quad \underline{0})$$



$$\theta_2 = (\underline{0.1} \quad \underline{0.1} \quad \underline{0.8})$$



$$\theta_3 = (\underline{0.3} \quad \underline{0.4} \quad \underline{0.3})$$

$$\begin{array}{|l} \rho_1 = (0.9 \quad 0.1) \\ \rho_2 = (0.1 \quad 0.9) \\ \rho_3 = (0.5 \quad 0.5) \end{array}$$

Example - simulation from two topics

Topic	Word distr.	probability	dna	gene	data	distribution
1	β_1	0.5	0.1	0.0	0.2	0.2
2	β_2	0.0	0.5	0.4	0.1	0.0

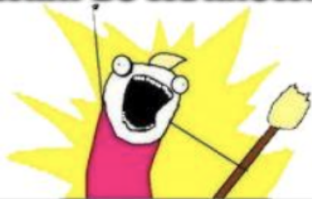
Doc 1	$\theta_1 = (0.2, 0.8)$				
	Word 1:	Topic=2	Word='gene'		
	Word 2:	Topic=2	Word='gene'		
	Word 3:	Topic=1	Word='data'		

Doc 2	$\theta_2 = (0.9, 0.1)$				
	Word 1:	Topic=1	Word='probability'		
	Word 2:	Topic=1	Word='data'		
	Word 3:	Topic=1	Word='probability'		

Doc 3	$\theta_2 = (0.5, 0.5)$				
⋮	⋮	⋮	⋮	⋮	⋮

Learning/inference in topic models

WHAT DO WE KNOW?



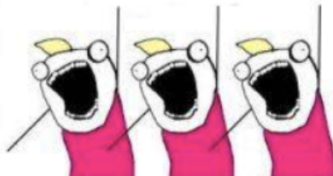
THE WORDS!



WHAT DO WE WANT?



**TOPICS PROPORTIONS, TOPIC
ASSIGNMENT AND WORD DISTRIBUTIONS!**



Learning/inference in topic models

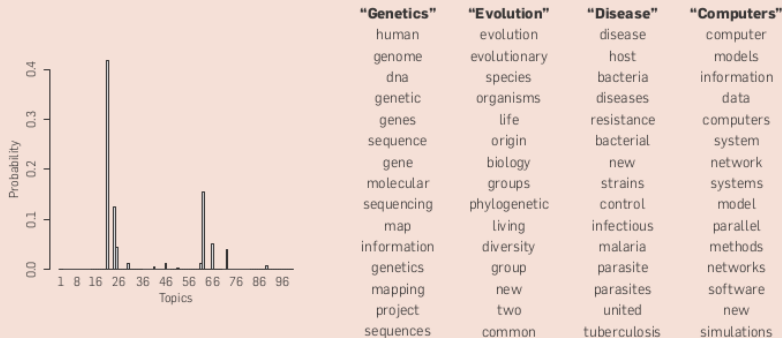
- What do we know?
 - ▶ The words in the documents: $w_{1:D}$
- What do we not know?
 - ▶ Topic proportions for each document: $\theta_{1:D}$
 - ▶ Topic assignments for each word in each document: $z_{1:D}$
 - ▶ Word distributions for each topic: $\beta_{1:K}$
- Do the Bayes dance: **Posterior distribution**

$$p(\theta_{1:D}, z_{1:D}, \beta_{1:K} | w_{1:D})$$

- The posterior is mathematically intractable:
 - ▶ **Gibbs sampling** (MCMC) [Correct, but can be slow]
 - ▶ **Variational Bayes** [approximate, but often fast]

Topic models Blei (2010)

Figure 2. Real inference with LDA. We fit a 100-topic LDA model to 17,000 articles from the journal *Science*. At left are the inferred topic proportions for the example article in Figure 1. At right are the top 15 most frequent words from the most frequent topics found in this article.



Gibbs sampler

- Bayes theorem

$$p(\theta|\text{data}) = \frac{p(\text{data}|\theta) \cdot p(\theta)}{p(\text{data})}$$

- For the topic model

$$\begin{aligned} p(z, \Theta, \Phi|w) &= \frac{p(z, \Theta, \Phi|w) \cdot p(z, \Theta, \Phi)}{p(w)} \\ &\propto p(z, \Theta, \Phi|w) \cdot p(z, \Theta, \Phi) \end{aligned}$$

Gibbs sampler

- Integrating out (collapsing) Θ and Φ :

$$p(z|w) = \int \int p(z, \Theta, \Phi|w) \cdot p(z, \Theta, \Phi) d\Phi d\Theta$$

will result in the following gibbs sampler

$$p(z_i = k | w_i, z_{-i}) = \underbrace{\frac{n_{k,v_i}^{(w)} + \beta}{n_{k,\cdot}^{(w)} + V\beta}}_{\text{type-topic } (\Phi)} \cdot \underbrace{(n_{k,d_i}^{(d)} + \alpha)}_{\text{topic-doc } (\Theta)}$$

where $n^{(w)}$ and $n^{(d)}$ are count matrices of size $D \times K$ and $K \times V$.

Example of $n^{(w)}$ and $n^{(d)}$

- Three ($D = 3$) documents:

w_1	boat	shore	bank
z_1	1	1	1

w_2	Zlatan	boat	shore	money	bank
z_2	2	1	1	3	3

w_3	money	bank	soccer	money
z_3	3	3	2	3

$$n^{(w)} = \begin{bmatrix} & \text{boat} & \text{shore} & \text{soccer} & \text{Zlatan} & \text{bank} & \text{money} \\ 2 & 2 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 2 & 2 \end{bmatrix}$$

$$n^{(d)} = \begin{bmatrix} 3 & 0 & 0 \\ 2 & 1 & 3 \\ 0 & 2 & 3 \end{bmatrix}$$

(Naive) algorithm

```
# Initialization
Sample all topic indicators randomly
Calculate  $\hat{n}(w)$  and  $\hat{n}(d)$ 

# Gibbs sampler
for each gibbs iteration do
  for each token  $w_i$  do
    remove  $z_i$  from  $\hat{n}(w)$  and  $\hat{n}(d)$ 
    for each  $k$  in 1 to  $K$  do
      
$$\text{prob}_k[k] = \frac{n_{k,v_i}^{(w)} + \beta}{n_{k,\cdot}^{(w)} + V\beta} \cdot (n_{k,d_i}^{(d)} + \alpha)$$

    end for
     $z_i \leftarrow \text{draw multinomial}(\text{prob}_k)$ 
    add  $z_i$  to  $\hat{n}(w)$  and  $\hat{n}(d)$ 
  end for
end for
return  $\hat{n}(w)$ ,  $\hat{n}(d)$ 
```

(Naive) algorithm

- Estimation of Φ and Θ

$$\hat{\phi}_{k,v} = \frac{n_{k,v}^{(w)} + \beta}{n_{k,\cdot}^{(w)} + V\beta}$$
$$\hat{\theta}_{d,k} = \frac{n_{d,k}^{(d)} + \alpha}{n_{d,\cdot}^{(d)} + K\alpha}$$

- Serial.
- Computational complexity is $O(K)$ for each token.
- Slow for larger corpuses.

Evaluation of topic models

- Convergence:

- ▶ **Log-likelihood**

- Evaluating and comparing models:

- ▶ **Held-out perplexity**

- ▶ See Wallach et al (2009). "Evaluation methods for topic models".