

Distributional Semantics

Lecture 3. Topic Modeling

Daniil Vodolazsky

March 9th, 2019

Lecture Plan

Lecture Plan

- 1 What Is Topic Modeling
- 2 Graphical Models
- 3 EM-algorithm
- 4 Probabilistic Generative Topic Model
- 5 Probabilistic LSA (PLSA)
- 6 Latent Dirichlet Allocation (LDA)
- 7 Additive Regularization of Topic Models (ARTM)
- 8 Seminar

What Is Topic Modeling

What Is Topic Modeling



President of Russia

Events

Structure

Videos and Photos

Documents

Contacts

Search

Meeting on the development of the Russia-Belarus cultural and humanitarian ties

Vladimir Putin and President of Belarus Alexander Lukashenko had a meeting at the Sirius Education Centre, where they discussed the prospects of developing bilateral cooperation in culture, education and sports.

February 15, 2019

14:00

Sochi

What is this text about? **politics** 60%, **education** 20%, **culture** 10%, **sports** 10%, **economics** 0%.

What Is Topic Modeling



What is the difference between topic modeling and classification?
Unspecified in advance topics!

Topic Modeling in Social Media

Актуальные темы: Россия

Изменить

Доренко

Сталина

Твитов: 8 256

Через 10

Твитов: 2 008

Дудя

На Кубани

Самарской

Твитов: 1 480

Международному

запад

Твитов: 1 413

В Кремле

Сергей

Твитов: 4 158

Актуальные темы: США

Изменить

#FatTuesday

Твитов: 6 886

#TuesdayThoughts

Твитов: 43 тыс.

King Kong Bundy

Wrestling legend King Kong Bundy has died, aged 61

#PancakeDay

Твитов: 39,4 тыс.

#TuesdayMotivation

Твитов: 18,9 тыс.

#PaczkiDay

Ty Cobb

Good Tuesday

Твитов: 11,1 тыс.

McCarthyism

Твитов: 25 тыс.

Former Trump White House

Твитов: 2 317

Актуальные темы: Франция

Изменить

#MardiGras

Твитов: 14,6 тыс.

#MardiConseil

Твитов: 1 883

#lesplanetes

Твитов: 4 012

Hugo Clément

Твитов: 2 220

#GILyon2019

Твитов: 1 504

Vieilles Charrues

Condé-sur-Sarthe

Твитов: 1 836

Sonic

Твитов: 114 тыс.

Black Eyed Peas

Твитов: 2 013

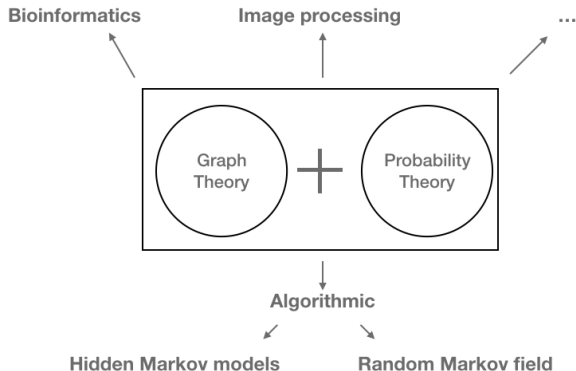
Bigard

Твитов: 11,4 тыс.

Graphical Models

Introduction to Graphical Models

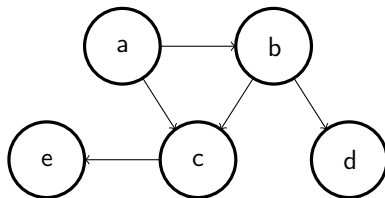
Graphical models bring together graph theory and probability theory in a powerful formalism for multivariate statistical modeling.



Directed Models

A **directed model** contains one factor for every random variable x_i in the distribution, and that factor consists of the conditional distribution over x_i given the parents of x_i , denoted $Pa_G(x_i)$:

$$p(\mathbf{x}) = \prod_i p(x_i \mid Pa_G(x_i)).$$



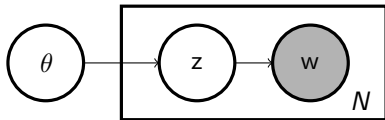
This graph corresponds to probability distributions that can be factored as follows:

$$p(a, b, c, d, e) = p(a)p(b \mid a)p(c \mid a, b)p(d \mid b)p(e \mid c).$$

Observable and Latent Variables

Observable variables are variables that can be observed and directly measured.

Latent (hidden) variables are variables that are not directly observed but are rather inferred from other variables that are observed (directly measured).



θ is a parameter, z_1, \dots, z_N are hidden (latent) variables, and w_1, \dots, w_N are observable variables.

Given θ , the joint probability for such model can be computed as follows:

$$p(\mathbf{w}, \mathbf{z} \mid \theta) = \prod_{i=1}^N p(z_i \mid \theta) p(w_i \mid z_i).$$

EM-algorithm

EM-algorithm is a powerful method that is widely used in probabilistic models. We will use \mathbf{X} to denote the observable variables, \mathbf{Z} to denote the latent variables, and Θ to denote the model parameters (that we want to estimate).

Assume that we know how to compute a joint distribution over \mathbf{X} and \mathbf{Z} given Θ : $p(\mathbf{X}, \mathbf{Z} \mid \Theta)$. Let's write a log-likelihood function as an expectation over \mathbf{Z} :

$$\log p(\mathbf{X} \mid \Theta) = \int q(\mathbf{Z}) \log p(\mathbf{X} \mid \Theta) d\mathbf{Z}$$

where $q(\mathbf{Z})$ is an arbitrary probability density for \mathbf{Z} .

$$\begin{aligned}\log p(\mathbf{X} | \Theta) &= \int q(\mathbf{Z}) \log p(\mathbf{X} | \Theta) d\mathbf{Z} = \int q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} | \Theta)}{p(\mathbf{Z} | \mathbf{X}, \Theta)} \frac{q(\mathbf{Z})}{q(\mathbf{Z})} d\mathbf{Z} \\ &= \int q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z} | \Theta)}{q(\mathbf{Z})} d\mathbf{Z} + \int q(\mathbf{Z}) \log \frac{q(\mathbf{Z})}{p(\mathbf{Z} | \mathbf{X}, \Theta)} d\mathbf{Z} \\ &= \mathcal{L}(q, \Theta) + D_{\text{KL}}(q(\mathbf{Z}) \| p(\mathbf{Z} | \mathbf{X}, \Theta)) \geq \mathcal{L}(q, \Theta).\end{aligned}$$

A full log-likelihood function is hard to optimize. Instead of maximizing $\log p(\mathbf{X} | \Theta)$ by Θ , we are going to maximize $\mathcal{L}(q, \Theta)$ by q, Θ . Notice that $\log p(\mathbf{X} | \Theta)$ does not depend on q so we are not limited to choose q . Thus we can turn the last inequality into equality by putting $q(\mathbf{Z}) = p(\mathbf{Z} | \mathbf{X}, \Theta)$ so that $D_{\text{KL}}(q(\mathbf{Z}) \| p(\mathbf{Z} | \mathbf{X}, \Theta)) = 0$.

- **E-step** creates a function for the expectation of the log-likelihood evaluated using the current estimate for the parameters:

$$q(\mathbf{Z})^{(n+1)} = \arg \max_{q'} \mathcal{L}(q', \Theta^{(n)}) = p(\mathbf{Z} \mid \mathbf{X}, \Theta^{(n)}).$$

- **M-step** computes parameters maximizing the expected log-likelihood found on the E-step:

$$\begin{aligned}\Theta^{(n+1)} &= \arg \max_{\Theta'} \mathcal{L}(q^{(n+1)}, \Theta') \\ &= \arg \max_{\Theta'} \int q(\mathbf{Z})^{(n+1)} \log \frac{p(\mathbf{X}, \mathbf{Z} \mid \Theta')}{q(\mathbf{Z})^{(n+1)}} d\mathbf{Z} \\ &= \arg \max_{\Theta'} \int q(\mathbf{Z})^{(n+1)} \log p(\mathbf{X}, \mathbf{Z} \mid \Theta') d\mathbf{Z} \\ &= \arg \max_{\Theta'} \mathbb{E}_{\mathbf{Z} \sim q(\mathbf{Z})^{(n+1)}} [\log p(\mathbf{X}, \mathbf{Z} \mid \Theta')].\end{aligned}$$

Probabilistic Generative Topic Model

Probabilistic Generative Topic Model

The joint probability of a document d and a word w can be expressed with conditional probability formula as

$$P(d, w) = P(d)P(w | d).$$

Now we will define a model with hidden variables $\mathcal{Z} = \{z_1, \dots, z_T\}$. We assume that each document has its own distribution over (latent) topics, and that each topic has its own distribution over (observable) words. So our goal is to find the matrices Φ and Θ such that

$$P(w | d) = \sum_{z \in \mathcal{Z}} P(w | z, d)P(z | d) = \sum_{z \in \mathcal{Z}} P(w | z)P(z | d) = \sum_{z \in \mathcal{Z}} \Phi_{w,z} \Theta_{z,d}.$$

In other words, we want to get a low-rank matrix factorization.

Probabilistic Generative Topic Model

A log-likelihood function for our model:

$$\begin{aligned}\log P(\mathcal{D}, \mathcal{W} \mid \Phi, \Theta) &= \log \prod_{d \in \mathcal{D}} \prod_{w \in \mathcal{W}} P(d, w)^{\#(w, d)} \\&= \sum_{d \in \mathcal{D}} \sum_{w \in \mathcal{W}} \#(w, d) \log P(d, w) \\&= \sum_{d \in \mathcal{D}} \sum_{w \in \mathcal{W}} \#(w, d) \log(P(d)P(w \mid d)) \\&= \sum_{d \in \mathcal{D}} \sum_{w \in \mathcal{W}} \#(w, d) \log \left(P(d) \sum_{z \in \mathcal{Z}} \Phi_{w,z} \Theta_{z,d} \right) \longrightarrow \max_{\Phi, \Theta},\end{aligned}$$

which is an objective function with restrictions:

$$\sum_{w \in \mathcal{W}} \Phi_{w,z} = 1, \Phi_{w,z} \geq 0, \sum_{z \in \mathcal{Z}} \Theta_{z,d} = 1, \Theta_{z,d} \geq 0.$$

Probabilistic Generative Topic Model

If $P(d)$ does not depend on Φ or Θ , then the problem is equivalent to

$$\sum_{d \in \mathcal{D}} \sum_{w \in \mathcal{W}} \#(w, d) \log \sum_{z \in \mathcal{Z}} \Phi_{w,z} \Theta_{z,d} \longrightarrow \max_{\Phi, \Theta}.$$

A point (Φ, Θ) of a local extremum can be found as a solution of the system:

$$\begin{aligned} P(z \mid d, w) &\propto \Phi_{w,z} \Theta_{z,d}, \\ \Phi_{w,z} &\propto \sum_{d \in \mathcal{D}} \#(w, d) P(z \mid d, w), \\ \Theta_{z,d} &\propto \sum_{w \in \mathcal{W}} \#(w, d) P(z \mid d, w). \end{aligned}$$

It means that on an E-step we should compute $P(z \mid d, w)$, and on an M-step we should compute $P(w \mid z) = \Phi_{w,z}$ and $P(z \mid d) = \Theta_{z,d}$ according to the formulas obtained.

PLSA

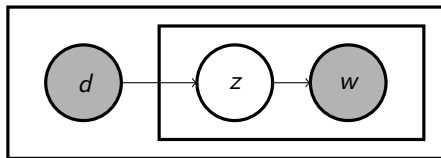
Probabilistic latent semantic analysis (PLSA) (1999) was the first technique in history for probabilistic topic modeling. The PLSA model is a latent variable model for co-occurrence data which associates an unobserved class variable $z \in \mathcal{Z} = \{z_1, \dots, z_T\}$ with each observation. A joint probability model over $\mathcal{D} \times \mathcal{W}$ is defined by the mixture

$$P(d, w) = P(d)P(w \mid d), \quad P(w \mid d) = \sum_{z \in \mathcal{Z}} P(w \mid z)P(z \mid d).$$

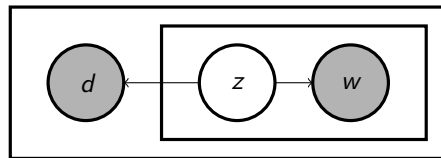
Since the cardinality of z is smaller than the number of documents/words in the collection, z acts as a bottleneck variable in predicting words. It is worth noticing that the model can be equivalently parameterized by:

$$P(d, w) = \sum_{z \in \mathcal{Z}} P(z)P(d | z)P(w | z)$$

which is perfectly symmetric in both entities, documents and words.



(a) non-symmetrical



(b) symmetrical

Two versions of the PLSA model.

PLSA model uses EM algorithm for maximum likelihood estimation.

- **E-step.** For each $z \in \mathcal{Z}$, $d \in \mathcal{D}$, $w \in \mathcal{W}$ compute

$$P(z \mid d, w) = \frac{P(d, z, w)}{P(d, w)} = \frac{P(z)P(d \mid z)P(w \mid z)}{\sum_{z' \in \mathcal{Z}} P(z')P(d \mid z')P(w \mid z')}.$$

- **M-step** For each $z \in \mathcal{Z}$, $d \in \mathcal{D}$, $w \in \mathcal{W}$ compute

$$P(w \mid z) \propto \sum_{d \in \mathcal{D}} \#(w, d)P(z \mid d, w),$$

$$P(d \mid z) \propto \sum_{w \in \mathcal{W}} \#(w, d)P(z \mid d, w),$$

$$P(z) \propto \sum_{d \in \mathcal{D}} \sum_{w \in \mathcal{W}} \#(w, d)P(z \mid d, w).$$

PLSA: Relation to LSA

Define matrices by $\hat{\mathbf{U}} : \hat{U}_{i,k} = P(w^{(i)} | z_k)$, $\hat{\mathbf{\Sigma}} : \hat{\Sigma}_{k,k} = P(z_k)$, $\hat{\mathbf{V}}^\top : \hat{V}_{k,j}^\top = P(d^{(j)} | z_k)$. The joint probability model \mathbf{P} can then be written as a matrix product $\mathbf{P} = \hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\hat{\mathbf{V}}^\top$.

The crucial difference between PLSA and LSA, however, is the objective function utilized to determine the optimal decomposition/approximation.

- In LSA, this is the L_2 or Frobenius norm, which corresponds to an implicit additive Gaussian noise assumption on counts.
- In contrast, PLSA relies on the likelihood function of multinomial sampling and aims at an explicit maximization of the predictive power of the model.

LDA

Latent Dirichlet allocation (LDA) (2003) is a topic modeling algorithm based on a generative graphical model. The concept of LDA is to consider all documents as composed of set of topics or clusters of words and to observe the distribution of them.

- We have a document collection \mathcal{D} and a vocabulary \mathcal{W} .
- Each document $d \in \mathcal{D}$ consists of N_d words $w_{d,1}, \dots, w_{d,N_d}$.
- The probability distributions over topics \mathcal{Z} are sampled for each document $d \in \mathcal{D}$ from a Dirichlet distribution with parameter α .
- For each topic $z \in \mathcal{Z}$ we generate a probability distribution over all words \mathcal{W} $\Phi_{:,z}$.
- After that, for each document d we generate N_d topics $z_{d,1}, \dots, z_{d,N_d}$ from distribution $\Theta_{:,d}$, and for each of these topics $z_{d,k}$ generate per a single word from distribution $\Phi_{:,z_{d,k}}$.

Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson at Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **simple numbers** game, particularly as more and more **genomes** are completely sequenced and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains

Arcady Mushegian, a **computational molecular biologist** at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

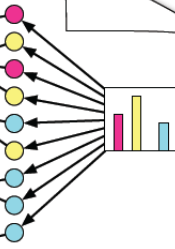


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments

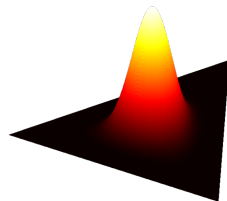
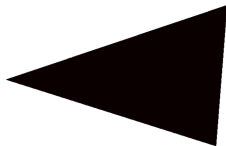
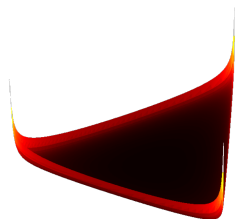


Dirichlet Distribution

A K -dimensional Dirichlet random variable θ can take values in the $(K - 1)$ -simplex (a K -vector θ lies in the $(K - 1)$ -simplex if $\theta_k \geq 0$ and $\sum_{k=1}^K \theta_k = 1$), and has the following probability density on this simplex:

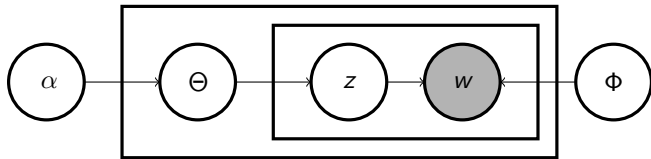
$$p(\theta \mid \alpha) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

where the parameter α is a K -vector with components $\alpha_k > 0$, and where $\Gamma(x)$ is the Gamma function.



Dirichlet distribution plots for $\alpha = (0.75, 0.75, 0.75), (1, 1, 1), (10, 10, 10)$

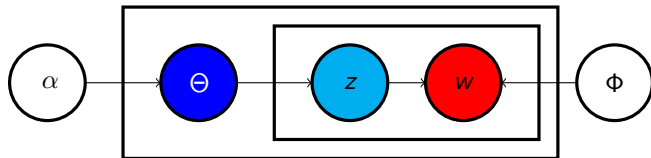
LDA Model



A graphical model for LDA.

In this model, α and Φ are parameters. Φ is a matrix of size $|\mathcal{W}| \times |\mathcal{Z}|$, Θ is a matrix of size $|\mathcal{Z}| \times |\mathcal{D}|$, α is a vector of size $|\mathcal{Z}|$. Moreover, $\Theta_{:,d} \sim \text{Dir}(\alpha)$.

LDA Model



A graphical model for LDA.

$$P(\mathcal{D}, \mathcal{W}, \mathcal{Z}, \Theta \mid \Phi, \alpha) = \underbrace{\prod_{d=1}^{|\mathcal{D}|}}_1 \underbrace{P(\Theta_{:,d} \mid \alpha)}_2 \underbrace{\prod_{n=1}^{N_d}}_3 \underbrace{P(z_{d,n} \mid \Theta_{:,d})}_4 \underbrace{P(w_{d,n} \mid z_{d,n}, \Phi)}_5 \longrightarrow \max_{\Phi, \alpha}$$

- 1. for each document
- 2. generate topic probabilities
- 3. for each word
- 4. select topic
- 5. select word from topic

$$\begin{aligned}
 P(\mathcal{D}, \mathcal{W}, \mathcal{Z}, \Theta \mid \Phi, \alpha) &= \prod_{d=1}^{|\mathcal{D}|} P(\Theta_{:,d} \mid \alpha) \prod_{n=1}^{N_d} P(z_{d,n} \mid \Theta_{:,d}) P(w_{d,n} \mid z_{d,n}, \Phi) \\
 &= \prod_{d=1}^{|\mathcal{D}|} \frac{\Gamma(\sum_{t=1}^{|\mathcal{Z}|} \alpha_t)}{\prod_{t=1}^{|\mathcal{Z}|} \Gamma(\alpha_t)} \prod_{t=1}^{|\mathcal{Z}|} \Theta_{t,d}^{\alpha_t-1} \prod_{n=1}^{N_d} \Theta_{z_{d,n},d} \Phi_{w_{d,n},z_{d,n}} \\
 &= \prod_{d=1}^{|\mathcal{D}|} \frac{\Gamma(\sum_{t=1}^{|\mathcal{Z}|} \alpha_t)}{\prod_{t=1}^{|\mathcal{Z}|} \Gamma(\alpha_t)} \prod_{t=1}^{|\mathcal{Z}|} \Theta_{t,d}^{\alpha_t-1} \prod_{n=1}^{N_d} \prod_{t=1}^{|\mathcal{Z}|} \Theta_{t,d}^{[z_{d,n}=t]} \Phi_{w_{d,n},t}^{[z_{d,n}=t]}.
 \end{aligned}$$

- **E-step.**

$$P(\Theta, \mathcal{Z} \mid \mathcal{D}, \mathcal{W}, \Phi, \alpha) \approx q(\Theta)q(\mathcal{Z}) = \arg \min_{q(\Theta), q(\mathcal{Z})} D_{\text{KL}}(q(\Theta)q(\mathcal{Z}) \parallel P(\Theta, \mathcal{Z} \mid \mathcal{D}, \mathcal{W}, \Phi, \alpha)).$$

- **M-step.**

$$\Phi = \arg \max_{\Phi} \mathbb{E}_{\Theta \sim q(\Theta), \mathcal{Z} \sim q(\mathcal{Z})} \log P(\Theta, \mathcal{Z} \mid \mathcal{D}, \mathcal{W}, \Phi, \alpha).$$

Questions?

Seminar

Topic Modeling on Hillary Clinton Emails

- 1 Reading data
- 2 Preprocessing
- 3 LSA
- 4 LDA
- 5 Homework: PLSA

<https://www.kaggle.com/s231644/topic-modeling-on-hillary-clinton-s-emails>

▶ Link