

Topic Modeling Library and Framework User Guide

Avanesov Valeriy
Kozlov Ilya

March 16, 2015

1 Input

The aim of this project is to provide flexible and simple library and framework for topic modeling. The main properties of this project:

- Ability to apply user-defined regularizer, sparsifier and stopping criteria.
- Ability to use multilingual topic modeling.

The paper is organized as follow: section 2 contains a quick introduction to topic modeling, 3 describes how to start a project without description of details of implementation. These two sections should allow you to use LDA and PLSA model in simple use-cases. The following section describes details of implementation and is necessary for deeper understanding of project structure and implementation of your own regularizer, sparsifier and other stuff. Section 4 describes a structure of the project, section 5 describes regularizers and explains how to implement your own regularizer, section 6 describes what sparsifier is and how to implement you own sparsifier. Section 7 for the description of how to define that convergence of EM-algorithm is achieved. Section 8 describes (actually, it doesn't, though we are working on it) how to use model in multilingual setting.

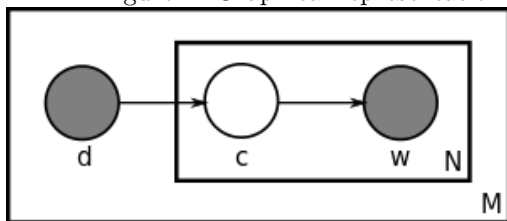
2 Introduction to topic modeling

2.1 Generative model

In topic modeling every document is viewed as a mixture of topics. Each topic is a multinomial distribution on words, so a generation model may be defined as follows:

- For every position in document d i.i.d choose topic t from distribution of topics by document
- Choose word w from topic t

Figure 1: Graphical representation of PLSA generative process



The aim of topic modeling is to recover topics and distribution of every document by topics.

2.1.1 Polylingual topic model

Suppose one has a collection of documents in different languages. We have a prior knowledge that some of the documents (or all of them) are written on the same topic but in different languages (for example one document may be a translation of another). Wikipedia may be considered as a source of this type of data. It leads us to a polylingual topic modeling [?]. In this model we assume that every topic is a set of multinomial distributions, one per language. Also we assume that every document may hold more than one set of words, so we represent document as a map from attributes \rightarrow sequences of words.

2.1.2 Robust PLSA

In robust PLSA we assume that some words are too specific for a document and can not be explained by topic distribution. Conversely, some of words are too common and may be explained by any topic. Robust PLSA take it into account. In robust PLSA a word may be generated from topic, noise or background. Noise is a multinomial distribution over rare words. Every document has its own noise. Background is a multinomial distribution over common words. Background is the same for every document. To generate words w in document d we:

- with probability $\propto \gamma$ generate word from noise

- with probability $\propto \varepsilon$ generate word from background
- with probability $\propto 1 - \gamma - \varepsilon$ generate word from topic distribution as in 2.1

γ and ε are hyperparameters of the model.

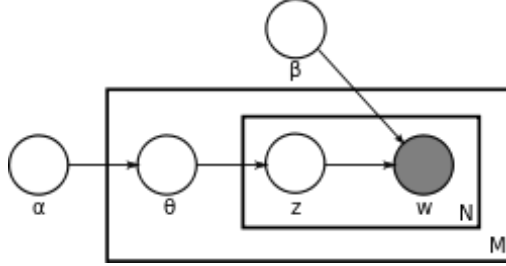
2.1.3 Sparse PLSA

A single document often corresponds to only a few topics, not to every. Analogously, word may corresponds only to a few topics, not to every topic. Sparse PLSA takes it into account, and replaces some of topics weights in document by zero and replaces some weights in distribution of words by topic by zero. Sparsification of topic modeling has some features which allow to sparse distribution of words by topic and distribution of document by topics without decreasing of model quality.

2.1.4 LDA

LDA is an extension of PLSA, LDA was described in [?] and now it is one of the most popular (may be the most popular) topic model. The idea of this model is that distribution of documents by topics and distribution of words by topics has Dirichlet¹ (usually, symmetric) prior distribution and we maximize marginal distribution $P(data, parameters)$ instead of conditional one - $P(data|parameters)$.

Figure 2: Graphical representation of LDA generative process



α and β are hyperparameters of the model.

2.2 Topic modeling as optimization problem

According to generative model one can estimate probability to observe collection D as:

$$p(D) = \prod_{d \in D} \prod_{w \in d} \sum_t p(t|d)p(w|t) \quad (1)$$

¹http://en.wikipedia.org/wiki/Dirichlet_distribution

Denote $\varphi_{wt} = p(w|t)$ and $\theta_{td} = p(t|d)$. One can obtain φ_{wt} and θ_{td} as a solution of optimization problem

$$L = \sum_{d \in D} \sum_{w \in d} \log \sum_t \varphi_{wt} \theta_{td} \rightarrow \max \quad (2)$$

with boundaries

$$\forall t \quad \sum_w \varphi_{wt} = 1, \quad \forall d \quad \sum_w \theta_{td} = 1 \quad (3)$$

and

$$\forall t, w \quad \varphi_{wt} \geq 0, \quad \forall d, t \quad \theta_{td} \geq 0 \quad (4)$$

2.3 Topic modeling as matrix decomposition

2.3.1 Kullback-Leibler divergence

Kullback-Leibler divergence is a non-negative semi-metric between a pair of probability distributions:

$$KL(p_i || q_i) = \sum_{i=1}^n p_i \ln \left(\frac{p_i}{q_i} \right) \quad (5)$$

Consider an empirical distribution \hat{p}_i and some parametric distribution $q_i = q_i(\alpha)$ which is used to explain \hat{p}_i . Easy to see that in this case minimization of KL-divergence is equivalent to estimation of α by maximum-likelihood:

$$KL(p_i || q_i(\alpha)) = \sum_{i=1}^n p_i \ln \left(\frac{p_i}{q_i(\alpha)} \right) \rightarrow \min_{\alpha} \Leftrightarrow \sum_{i=1}^n p_i \ln(q_i(\alpha)) \rightarrow \max_{\alpha} \quad (6)$$

Thus one can easily see that (2) equivalent to weighted Kullback-Leibler divergence minimization:

$$\sum_{d \in D} n_d KL_w \left(\frac{n_{dw}}{n_d} || \sum_{t \in T} \varphi_{wt} \theta_{td} \right) \rightarrow \min_{\Phi, \Theta} \quad (7)$$

where n_{wd} – number of words w in document d , n_d – number of words in document d .

2.3.2 Matrix decomposition

Denote empirical distribution of words by document as $\hat{p}(w, d) = \frac{n_{wd}}{n_d}$. According to this notation one can consider the problem (2) as matrix decomposition:

$$F \approx_{KL} \Phi \Theta \quad (8)$$

where matrix $F = (\hat{p}(w, d))_{W \times D}$ is empirical distribution of words by document, matrix $\Phi = (\varphi_{wt})_{W \times T}$ is distribution of words by topics and matrix

$\Theta = (\theta_{td})_{T \times D}$ is distribution of documents by topics. Thus, our optimization problem may be rewritten in Kullback–Leibler notation as

$$KL(F, \Phi\Theta) \rightarrow \min \quad (9)$$

Thus PLSA may be observed as stochastic matrix decomposition.

2.4 Expectation-Maximization algorithm

Unfortunately (2) has no analytical solution. Thus we use Expectation - Maximization (EM) algorithm. This algorithm consists of two steps:

1. Estimation of number n_{dwt} of words w , produced by topic t in document d . (E - step)
2. Optimization of distribution of documents by topics and optimization of distribution of topics by words relying on the n_{dwt} values obtained during E - step . (M - step)

One can estimate n_{dwt} as follows:

$$n_{dwt} = n_{wd}p(w|t)p(t|d) \quad (10)$$

where n_{wd} – number of words w in document d . Thus, probability $p(w|t)$ may be estimated as

$$p(w|t) = \frac{n_{wt}}{n_t} = \frac{\sum_d n_{dwt}}{\sum_w \sum_d n_{dwt}} \quad (11)$$

Similarly for $p(t|d)$

2.4.1 Complexity

Obviously, time complexity for every iteration is $O(T \times n_{uniq} \times D)$.

Space complexity is given by $O(W \times T + D \times T)$.

Here n_{uniq} is the average number of distinct tokens in a document, D stands for number of documents and W – for the size of vocabulary.

2.5 Regularizers

Regularizers may improve human-understandability of the topics, transform PLSA to LDA, provide an ability for semi-supervised learning (employ a prior knowledge about document topic or topics structure), select number of topics. Instead of optimization (2) we optimize

$$L(\Phi, \Theta) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta} \quad (12)$$

Where $R(\Phi, \Theta)$ is a twice differentiable function, named regularizer. Solution of this problem leads to a modification of M-step:

$$\varphi_{wt} \propto \left(\hat{n}_{wt} + \varphi_{wt} \frac{\partial R(\Phi, \Theta)}{\partial \varphi_{wt}} \right)_+, \quad \theta_{td} \propto \left(\hat{n}_{dt} + \theta_{td} \frac{\partial R(\Phi, \Theta)}{\partial \theta_{td}} \right)_+ \quad (13)$$

Where \hat{n} has been estimated by E-step.

2.5.1 Probability interpretation

Does regularizer have probabilistic interpretation? Yes, it does. Imagine that Φ and Θ have some prior distribution, for example dirichlet distribution. In this case:

$$P(D, \Theta, \Phi) = P(D|\Theta, \Phi) \times P(\Theta, \Phi) \quad (14)$$

Take the logarithm of the equation 14:

$$\ln P(D, \Theta, \Phi) = \ln(P(D|\Theta, \Phi) \times P(\Theta, \Phi)) = \ln(P(D|\Theta, \Phi)) + \ln(P(\Theta, \Phi)) \quad (15)$$

Denote $L(\Phi, \Theta) = \ln(P(D|\Theta, \Phi))$ and $R(\Phi, \Theta) = \ln(P(\Theta, \Phi))$ and obtain 12

Example Consider the example with symmetric dirichlet distribution² with parameter α . Assume we estimate n_{td} for each topic in E-step and want to perform M- step. We would find new distribution of document by topic as a maximum posterior probability solution:

$$\ln(P(\vec{\theta}|\vec{n}_{dt})) \rightarrow \max \quad (16)$$

$$\ln(P(\vec{\theta}|\vec{n}_{dt})) \propto \ln(P(\vec{n}_{dt}|\vec{\theta}) \times P(\vec{\theta})) \Leftrightarrow \quad (17)$$

$$\ln(P(\vec{\theta}|\vec{n}_{dt})) = \ln\left(\prod_{t \in T} \theta_t^{n_{td}} \times \prod_{t \in T} \theta_t^{\alpha-1}\right) + \text{const} \Leftrightarrow \quad (18)$$

$$\ln(P(\vec{\theta}|\vec{n}_{dt})) = \left(\sum_{t \in T} (n_{td} + \alpha - 1) \ln(\theta_t)\right) + \text{const} \Leftrightarrow \quad (19)$$

And boundary:

$$\sum_{t \in T} \theta_t = 1 \quad (20)$$

We may write down a Lagrangian from 19 and 20

$$L(\vec{\theta}, \lambda) = \sum_{t \in T} (n_{td} + \alpha - 1) \ln(\theta_t) - \lambda \left(\sum_{t \in T} \theta_t - 1\right) \quad (21)$$

$$\frac{\partial L}{\partial \theta_t} = \frac{n_{td} + \alpha - 1}{\theta_t} - \lambda = 0 \quad (22)$$

$$\theta_t \propto n_{td} + \alpha - 1 \quad (23)$$

Where n_{td} was estimated in E-step and $\theta_{td} \frac{\partial R(\Phi, \Theta)}{\partial \theta_{td}} = \theta_{td} \frac{\partial (\alpha-1) \ln(\theta_{td})}{\partial \theta_{td}} = \alpha - 1$ is regularizer.

See example of regularizer implementation in `ru.ispras.modis.tm.regularizer.SymmetricDirichlet`

²http://en.wikipedia.org/wiki/Dirichlet_distribution

3 Quick start

In order to use topic modeling one have to perform the following step:

1. Read documents
2. Split each document into a sequence of words
3. Replace words by it serial number
4. Build a topic model.
5. Train the topic model.
6. Save results or use it in application.

First of all one may look at `tm/src/main/scala/ru/ispras/modis/tm/scripts/QuickStart.scala`
In this section would be shown how to perform each step.

3.1 Read the data

First of all we have to read the data from disc. It worth mentioning that text normalization is non-goal of our project, thus text should be already preprocessed. This step depends on the organization of input data, so it is no use to describe this step in general way. Instead of this we describe this step for our example file, one may easy modify this step for his own use case. In this introduction we would use file `arxiv.part` from directory `examples`. This is a thousand of scientific papers from Arxiv³. Texts was preprocessed, words were separated by space. Every line corresponds to a single document.

First of all one has to read data, and split each line by space:

```
val lines = Source.fromFile(new File("examples/arxiv.part")).getLines()
val wordsSequence = lines.map(line => line.split(" "))
```

Than one have to construct a `TextualDocument` from each sequence of words:

```
val textualDocuments = wordsSequence.map(words =>
  new TextualDocument(Map(DefaultAttributeType -> words)))
```

Each `TextualDocument` is a map from attribute to the correspond sequence of words. For example attribute may denote language of text in the case of multilingual topic modeling. If there is only one attribute (as in our case) one can use `DefaultAttributeType`. Furthermore, in such a case one can use `SingleAttributeNumerator` instead on `Numerator` – it takes `Seq[String]` as inputs and implicitly uses `DefaultAttributeType`.

3.2 Replace words by it serial number

In order to save the memory in our application we have to replace words by it serial number and group the same words together.

$$Seq(duck, duck, duck, goose) \rightarrow Seq((0,3),(1,1))$$

³<http://arxiv.org/>

For this purpose we may use object `Numerator`. It takes into input sequence of textual documents, replace words by it serial number, group the same words and return sequence of documents. It also return `Alphabet` (map from word to serial number and vice versa) `val (documents, alphabet) = Numerator(textualDocuments).`

3.2.1 How to use alphabet

`Alphabet` holds a map from words to their indices and vice versa. Thus it allow to

1. get words by it index and attribute:
`alphabet(Category, 1) // goose`
2. get index of word by attribute and word:
`alphabet.getIndex(Category, duck) // 0`
3. get the number of words, corresponding to the given attribute
`alphabet.numberOfWords(Category) // 100500`

3.3 Build model

One can build model with class `Builder` (see package `ru.ispras.modis.builder`). Our project provides three types of builders:

1. `PLSABuilder` – builds a standard PLSA (as it is described in original paper [?])
2. `LDABuilder` – builds an LDA (PLSA with Dirichlet regularizer, for details section 2.5)
3. `RobustPLSABuilder` – builds a robust PLSA. Robust PLSA takes into account that some words are too rare and can't be explained by any topic (we call this kind of words "noise"). Some other words may be too common (for example a stop-word, that we forget to remove in the preprocessing step). Words of this kind are referred to any topic (we call that kind of word "background").

In this example we would use a standard PLSA, other builders work analogously. To build `plsa` one should

1. Set number of topics in model:
`val numberOfTopics = 25`
2. Set the number of iterations for EM-algorithm to perform:
`val numberOfIteration = 100`
3. Create an instance of class `java.util.Random`:
`val random = new Random()`

4. And build the model:

```
val builder = new PLSABuilder(numberOfTopics, alphabet, documents,
    random, numberOfIteration)
```

3.4 Training a model

Now we build a model and can perform stochastic matrix decomposition(train the model)

```
val trainedModel = plsa.train(documents)
```

trainedModel holds the matrices Φ and Θ (see 2.3) Φ is distribution of words by topic thus the number in the intersects of i – th row and j – th column show the probability to generate word j from topic i . Each attribute set according to the matrix Φ . To obtain matrix Φ , corresponds to attribute Category:

```
val phi = trainedModel.phi(Category)
```

Θ is a distribution of documents by topic thus the number in the intersects of i – th row and j – th column show the weight of topic j in document i . To obtain matrix Θ :

```
val theta = trainedModel.theta
```

3.4.1 How to interpret the results

If you are not familiar with topic modeling read 2.1. In our library we store the distribution of words by topic in matrix Φ , distribution of documents by topics in matrix Θ . In order to obtain the probability to generate word w from topic t one may use method

```
trainedModel.phi(Category).probability(t, w)
```

In order to obtain the probability to generate topic t in document d

```
trainedModel.theta.probability(d, t)
```

3.4.2 How to deal with matrices

Matrices Φ and Θ hold different information, but support the similar methods. Both classes hold matrix of expectation and stochastic matrix. Expectation matrix holds values, estimated by E-step (see 2.4), stochastic matrix contain a probability. Method of classes Phi and Theta are explained in table 1

We also have an utility to save matrix in file:

```
TopicHelper.saveMatrix("path/to/file", matrix) where matrix may be  $\Phi$  or  $\Theta$ . One also may write top  $n$  words from each topic to estimate topic coherence with TopicHelper.printAllTopics(n, phi, alphabet)
```

Table 1: Method of class Phi and class Theta

method	Phi	Theta
addToExpectation(<i>row</i> , <i>column</i> , <i>value</i>)	<i>row</i> - topic index column - word index <i>column</i> - topic index	<i>row</i> - document index value - n_{dwt} value - n_{dwt}
probability(<i>rowIndex</i> , <i>columnIndex</i>)	probability to generate word <i>columnIndex</i> from topic <i>rowIndex</i>	weight of topic <i>columnIndex</i> in document <i>rowIndex</i>
expectation(<i>rowIndex</i> , <i>columnIndex</i>)	return expected value for words, <i>columnIndex</i> generated from topic <i>rowIndex</i> : n_{wt}	return expected value for words, generated from topic <i>columnIndex</i> in document <i>rowIndex</i> : n_{dt}
numberOfRows	number of topics	number of documents
numberOfColumns	number of words	number of topics
dump()	replace negative value in expectation matrix by zero, replace value in stochastic matrix by corresponding values from expectation matrix, normalize stochastic matrix and replace values in expectation matrix by zero	

4 Building bricks of our model

In the previous section we have described how to build a model, but we still have not described how our implementation works. Let's fix this shortcoming. Our model supports a multilingual documents, thus document may contain few texts related to different attributes, every text from the same document inherits the same distribution of topics. Thus we have one matrix Φ per attribute and one matrix Θ for all attributes. In order to train our model we have to

1. Generate some initial approximation for matrix Θ and every matrix Φ
2. Perform E-step and estimate the number of words w in document d , produced by topic t for every attribute.
3. Apply regularizer to every matrix. (see 2.5)
4. Perform M-step for every matrix.
5. Sparsify matrices Φ and matrix Θ . (see 2.1.3)
6. Check the stopping criteria and return training model if it time to stop or return to the step 2 otherwise.

The main part of our project is PLSABricks, it does the main part of work. Every brick processes a single attribute. It performs E-step, applies regularizer to matrix Φ , performs M-step to the matrix Φ and sparsifies matrix Φ . PLSA includes one brick per attribute. Also it stores regularizer and sparsifier for matrix Θ and stopping criteria.

5 Regularizer

5.1 Regularizer

The theory is given in section 2.5. Here we describe the details of implementation of regularizer in our project and describe some standard regularizers, which are implemented in our library.

5.1.1 Implementation

One can find regularizers in `ru.ispras.modis.regularizer`. Any regularizer should inherit from class `Regularizer`. In order to implement his own regularizer user have to implement methods `regularizePhi`, `regularizeTheta` and `apply`.

Implementation of `regularizePhi`

As one may see in 13,

$$\varphi_{wt} \propto \left(\hat{n}_{wt} + \varphi_{wt} \frac{\partial R(\Phi, \Theta)}{\partial \varphi_{wt}} \right)_+ \quad (24)$$

Thus, in order to implement method `regularizePhi` one have to:

- calculate $\varphi_{wt} \frac{\partial R(\Phi, \Theta)}{\partial \varphi_{wt}}$ using matrix Φ and matrix Θ for each word w and each topic t .

In order to take i -th row and j -th column in matrix Φ one may use

$$phi.probability(i, j)$$

where i – topic index, j – word index.

Analogously for matrix Θ :

$$theta.probability(i, j)$$

, where i -topic index, j - document index.

- Add these values to matrix of expectation, in order to add $\varphi_{wt} \frac{\partial R(\Phi, \Theta)}{\partial \varphi_{wt}}$ to expectation matrix n_{wt} one should use method

$$phi.addToExpectation(t, w, \varphi_{wt} \frac{\partial R(\Phi, \Theta)}{\partial \varphi_{wt}})$$

Implementation of `regularizeTheta`

This method is analogous to the previous paragraph. One have to calculate

$$\theta_{td} \frac{\partial R(\Phi, \Theta)}{\partial \varphi_{td}}$$

and add it to expectation of n_{dt} .

Implementation of apply

apply is used to calculate ℓ_2 instead of log likelihood (and corresponding perplexity). If you want calculate log likelihood or you are lazy to implement this method return 0f.

5.1.2 Implemented regularizer

Now in our project are implemented following regularizers:

- ZeroRegularizer, it regularizer do nothing, if you don't want to use regularizer use this one
- RegularizerSum allow to apply a sequence of regularizers sequentially. For example if you have a few regularizer: r_1, r_2, r_3 and you want to apply them sequentially. For this aim

```
import ru.ispras.modis.regularizer.Regularizer.toRegularizerSum  
val regularizerSum = r1 + r2 + r3
```
- SymmetricDirichlet, it regularizer add a dirichlet prior to the distribution of document by topic and words by topic, it is used to convert PLSA into LDA (see 2.1.4)

6 Sparsifier

One document often corresponds to only a few topics, not to everyone. Analogously, word corresponds only to a few topics, not to every topic. Thus some weights in matrices Φ and Θ may be replaced by zero without drop of quality of our model. For this purpose we use sparsifier. Any sparsifier should inherit from trait `Sparsifier` and implement a method `apply`. This method takes into input `MatrixForSparsifier` and `numberOfIteration`. Sparsifier may obtain probabilities from matrix and decide which cells should be replaced by zero, based on value, number of iteration and its state. To replace i -th row and j -th column by zero one should call `matrix.setZero(i, j)`

6.1 Implemented sparsifier

We implement

- `ZeroSparsifier`, it does nothing. If you don't want to use sparsifier you should use this one.
- `ThresholdSparsifier` - this class replaces value with zero if it is less than threshold if it did not replaced in this way more than `maxNumberOfZeroised` during this method call and the current number of iteration exceeds `startIteration`.

7 Stopping Criteria

There is no conventional way to define how many iteration should be performed. One of the common way to do that is to execute a fixed number of iterations, where the number of iteration is a some heuristic. There are other stopping criterion, thus we allow user to implement his own criterion. Any criterion should inherit criterion `trait StoppingCriteria`. Stopping criterion decides to stop based on perplexity after this iteration, perplexity after previous iteration, number of iterations and internal state.

8 Multilingual topic modeling

References

- [1] David Mimno, Hanna M. Wallach, Jason Naradowsky, David A. Smith, and Andrew McCallum. Polylingual topic models. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 2 - Volume 2*, EMNLP '09, pages 880–889, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics.
- [2] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, March 2003.
- [3] Thomas Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22Nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '99, pages 50–57, New York, NY, USA, 1999. ACM.