# Quick Guide

## 1. Loading and preprocessing of data.

Run **TopicMiner.exe** file (in TopicMiner directory). Press the button with a file icon to load textual data.



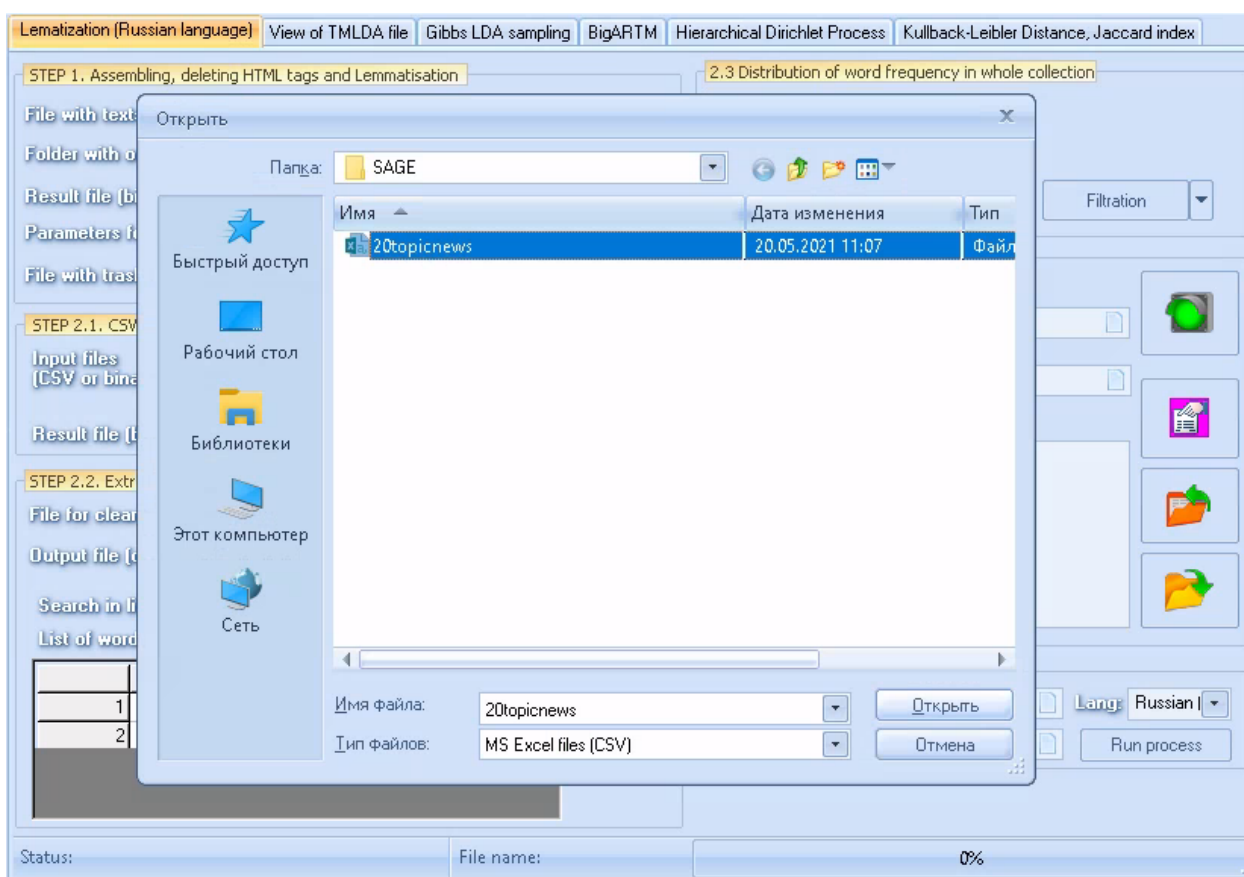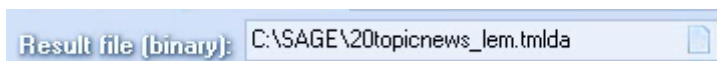Next, you will need to select a file with textual data, in which one line is one document.

```
1   docs
2   cs okstate chong kermit available windows article steve frampton
    wondering kermit package actual package usual ftp sites chong
3   usc bin looking address noise cancellation tech am new newsgroup
    ask question am looking address noise rather important help me
    regard please thank aludra usc
4   bear tigger cs colorado bear giles secret source molitor amolitor
    nmsu wrote monitor phonecalls monitor usenet may collect data
    making sense another matter sci crypt m graduate cs major strong
    math background programmer taking cryptology course
5   pwb aerg canberra au paul blackman moving article rutgers viamar
    kmembry remember reading program made windows icons run awayfrom
    mouse moved near them does anyone know nameof program ftp location
    probably cica file zip ll find icons cica line description your
    icons o paul blackman pwb science canberra au o water research
    centre pwb aerg canberra au o faculty applied science o university
    canberra australia spend little love get high lenny kravitz
6   cohen ssdgwy mdc andy cohen single launch space article mcimail
    karl dishaw wrote andy cohen single launch core station concept
    shuttle external tank solid rocket boosters used launch station
    into orbit shuttle main engines mounted tail station module launch
    jettisoned after et separation why jettison ssmes why hold them
```

For example, open a CSV file (see picture).



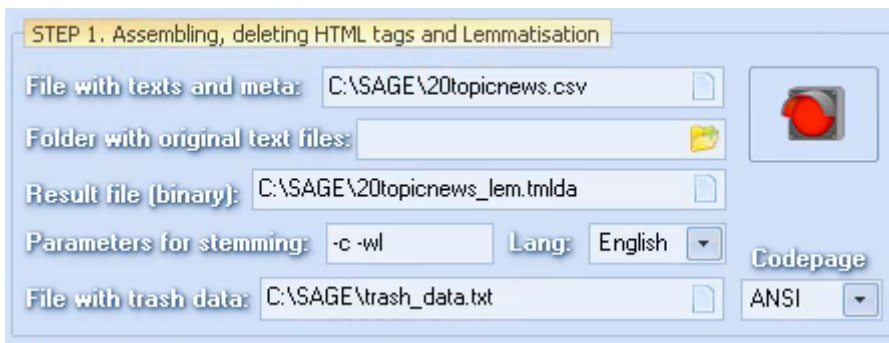Next, you need to specify the file to store the lemmatization results.



For example, **C:\SAGE\20topicnews_lem.tmlda**. Next, you need to specify text encoding (in this version, two options are implemented). Select **ANSI**.
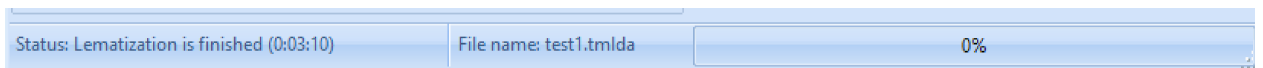


You also need to specify the language (Russian and English are implemented in this version). Select **English**.

Lastly, you need to specify a file that contains some trash words/entities, which can also be deleted. See **trash_data.txt** file for an example of such a file.

As a result, when everything you need is filled in, click on the red button.

You will start the process of lemmatization (see the percentage of completion at the bottom of the screen). As a result, you will have a file containing the original and lemmatized texts. When the first stage of preprocessing (lemmatization procedure) is over, you will see the following message in the status bar.
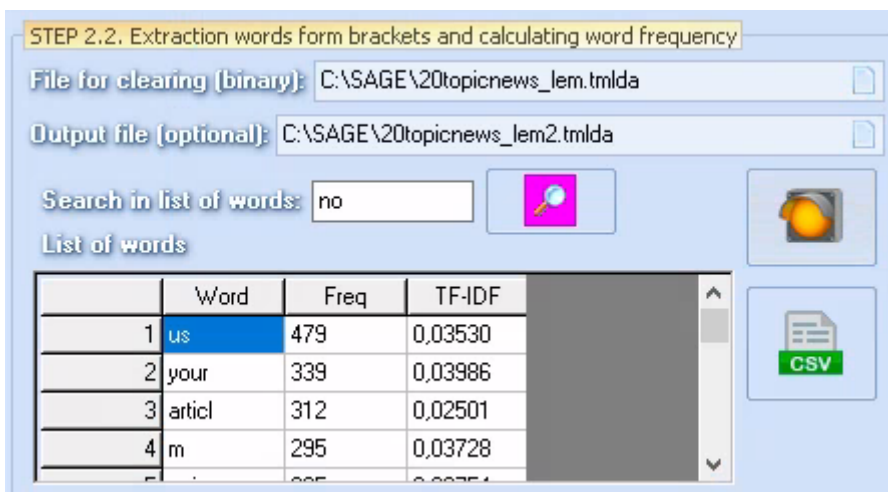


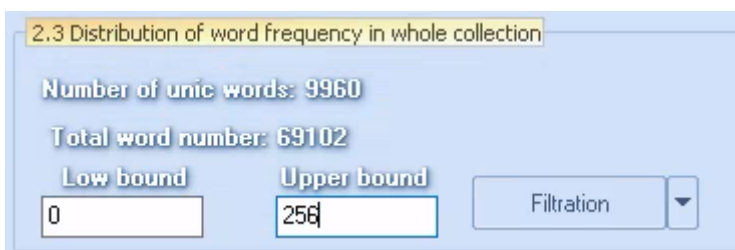The time of lemmatization is also indicated there.

At the second stage, you need to specify **C:\SAGE\20topicnews_lem.tmlda** as the input file, and **C:\SAGE\20topicnews_lem2.tmlda** as the output (it is optional). Then press the yellow button.



The process of the second stage (in percentage) is also visualized. A word list with frequencies will also appear. This list can be saved by pressing "CSV" button. This option is useful to form a list of stop words. For example, you can select the most frequent words as a list of stop words and put a filter that will save these words to a test file.



For example, in this case, let's take **256** as the upper bound.



Next, click on the "filtration" button and specify the file in which the words that lie outside the specified range will be saved (that is, everything that is above the frequency of **1258**). If you open such a file, you will see a list of stop words (most frequent).

At the last stage of preprocessing. You need to point to the input file
**C:\SAGE\20topicnews_lem2.tmlda**, and the output to
**C:\SAGE\20topicnews_lem3.tmlda**. And also download a file with a list of stop
words.



Now press the green button.

## 2. **Preprocessing results.**

If you need to see what you have learned (and also to start filtering texts), then you need to go to the next tab. To observe the results, use the following button:



Specify **C:\SAGE\20topicnews_lem3.tmlda** file.



Table will be displayed:

The first column is the original data, the second column is the lemmatized data. The third and subsequent columns are metadata.

As a result of preprocessing, empty documents can be obtained (for example, due to unspecified line ends or words in the texts were deleted using stop words). Of course, such documents must be removed from the collection. This can be done using filtering.

For example, let's delete short documents that contain only one word (**00000**):

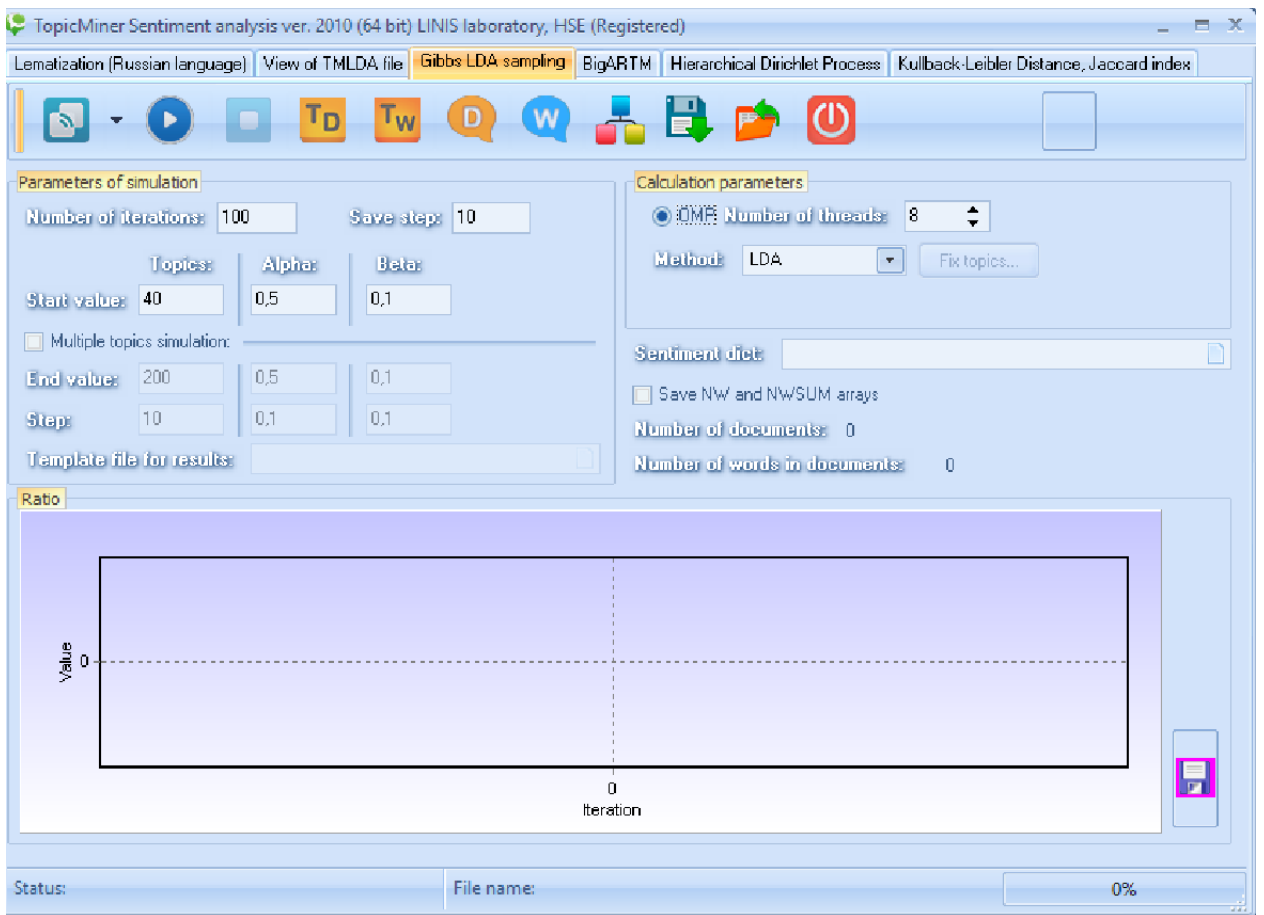| Document ID | Original document | Lematized document | Author | Field 1 |
|---|---|---|---|---|
| 445 | pierson enet dec dave pierson swr meter cb article peter m insane apana org au peter tryndoch allthe devil meter cb | pierson enet dec dave pierson swr meter cb articl peter m insan apana org au peter tryndoch allth devil meter cb radio | 446 | |
| 446 | gardner convex steve gardner escrow database article strnlght netcom david sternlight after waco massacre big | gardner convex steve gardner escrow databas articl strnlght netcom david sternlight after waco massacr big brother | 447 | |
| 447 | ken sugra uucp kenneth ng hst servicing mission scheduled daysin article hathaway stsci also implied other posters | ken sugra uucp kenneth ng hst servic mission schedul daysin articl hathawai stsci also impli other poster why need | 448 | |
| 448 | {00000} | 00000 | 449 | |
| 449 | jhan debra dgbt doc ca jerry han overreacting once tapped your code good any more article steve b access | jhan debra dgbt doc ca jerri han overreact onc tap your code good ani more articl steve b access digex steve | 450 | |
| 450 | mart csri toronto mart changing oil self bobml mxmsd msd measurex bob lagesse long silly discussion deleted while why | mart csri toronto mart chang oil self bobml mxmsd msd measurex bob lagess long silli discuss delet while why bother | 451 | |
| 451 | chin ee ualberta ca chin need info dsp want start dsp project music stereo cassette any chip set development kit | chin ee ualberta ca chin need info dsp want start dsp project music stereo cassett ani chip set develop kit compil | 452 | |
| 452 | darice yoyo cc monash au fred rice slavery why sex only allowed marriage guncer enuxha eas asu selim guncer | daric yoyo cc monash au fred rice slaveri why sex onli allow marriag guncer enuxha ea asu selim guncer might like | 453 | |
| 453 | mnhcc cunyvm bitnet marty helgesen public private revelation formerly question virgin ashley account private | mnhcc cunyvm bitnet marti helgesen public privat revel formerli question virgin ashlei account privat revel doe some | 454 | |

To do this, we will indicate the number of words:



And click on **Exclude short docs** button. As a result, **20topicnews_lem3_we.tmlda** file will be created, which will no longer contain these documents.
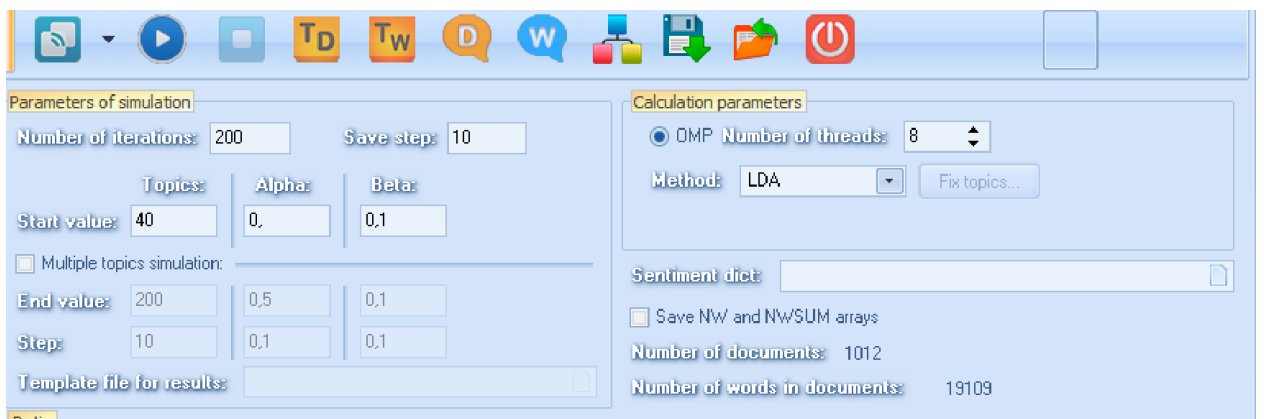
## 3. Topic modeling.

In order to start thematic modeling, for example, based on Gibbs sampling (Monte Carlo method), you need to go to another tab:

You need to load the file after preprocessing and removing empty documents.

For example, load the file **C:\SAGE\20topicnews_lem3_we.tmlda**. Ignore the download percentage (in the current version, this is just a measure of the length of the collection).

Next, specify the parameters of the model (number of topics, number of iterations, the rendering step, and *alpha*, *beta* parameters). This can be done as shown in the picture.
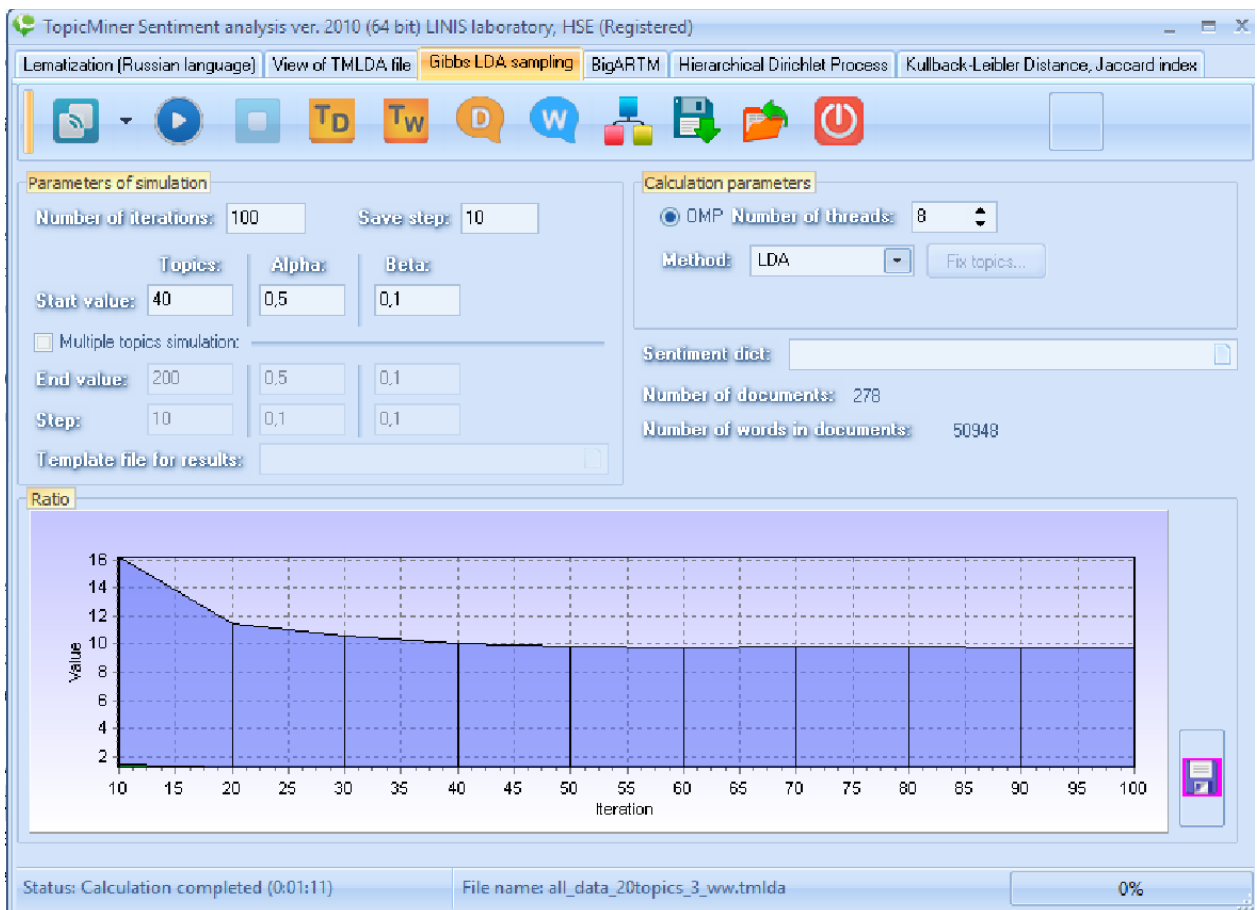


And you need to choose a model, for example, **LDA**.

Then click on the button:



The modeling progress will be displayed on the graph:



The blue line is the percentage of words with high probability. It can be seen from the graph that after 100 iterations, the percentage of words stops changing, which means a further increase in iterations does not need to be done. The green line is the percentage of high probability documents.

## 4. Viewing the results obtained and uploading to an external file.

To see the already sorted calculation results, you need to click on the button:

As a result, you will see the following picture:



If you want to read the most probable documents in any topic, then you can click on the button:



You will see the following:

Each cell contains the document ID and its probability in the topic (that is, in the column). To view the contents of a document (original non-lematized text), just place the mouse cursor onto the cell.

## 5. Saving calculation results as a project:

The calculation results can be saved as a project. The next time, this project can be opened in TopicMiner and you can continue working with it. To save the project, you need to use the following button:



And to open the project use the following button: