

# Topic Modelling

Topic modelling with Latent Dirichlet Allocation, Latent Semantic Indexing or Hierarchical Dirichlet Process.

## Inputs

- Corpus: A collection of documents.

## Outputs

- Corpus: Corpus with topic weights appended.
- Topics: Selected topics with word weights.
- All Topics: Token weights per topic.

**Topic Modelling** discovers abstract topics in a corpus based on clusters of words found in each document and their respective frequency. A document typically contains multiple topics in different proportions, thus the widget also reports on the topic weight per document.

The widget wraps gensim's topic models (**LSI**, **LDA**, **HDP**).

The first, LSI, can return both positive and negative words (words that are in a topic and those that aren't) and concurrently topic weights, that can be positive or negative. As stated by the main gensim's developer, Radim Řehůřek: *"LSI topics are not supposed to make sense; since LSI allows negative numbers, it boils down to delicate cancellations between topics and there's no straightforward way to interpret a topic."*

LDA can be more easily interpreted, but is slower than LSI. HDP has many parameters - the parameter that corresponds to the number of topics is *Top level truncation level (T)*. The smallest number of topics that one can retrieve is 10.

Topic	Topic keywords
1	came, little, went, one, king, go, away, father, man, took
2	little, hansel, king, gretel, children, forest, red, mother, woman, princess
3	bird, hansel, mother, tree, gretel, kywitt, juniper, beautiful, sing, wife
4	wife, shudder, man, youth, fish, fisherman, princess, fox, fire, learn
5	wife, fish, fisherman, shudder, youth, father, home, fire, go, bird
6	red, hansel, gretel, wolf, rose, bear, cap, snow, white, grandmother
7	fox, bird, pick, wolf, ashputtel, bride, tree, mother, golden, prince
8	falada, maid, curdken, head, blow, bride, fox, alas, sparrow, sadly
9	cap, grandmother, wolf, bear, white, snow, rose, little, children, dwarf
10	hedgehog, hare, field, cap, grandmother, run, wife, water, king, little

1. Topic modelling algorithm:
  - **Latent Semantic Indexing**. Returns both negative and positive words and topic weights.
  - **Latent Dirichlet Allocation**
  - **Hierarchical Dirichlet Process**
2. Parameters for the algorithm. LSI and LDA accept only the number of topics modelled, with the default set to 10. HDP, however, has more parameters. As this algorithm is computationally very demanding, we recommend you to try it on a subset or set all the required parameters in advance and only then run the algorithm (connect the input to the widget).
  - First level concentration ( $\gamma$ ): distribution at the first (corpus) level of Dirichlet Process
  - Second level concentration ( $\alpha$ ): distribution at the second (document) level of Dirichlet Process
  - The topic Dirichlet ( $\alpha$ ): concentration parameter used for the topic draws
  - Top level truncation ( $T$ ): corpus-level truncation (no of topics)
  - Second level truncation ( $K$ ): document-level truncation (no of topics)
  - Learning rate ( $\kappa$ ): step size
  - Slow down parameter ( $\tau$ )
3. Produce a report.
4. If *Commit Automatically* is on, changes are communicated automatically. Alternatively press *Commit*.

## Exploring Individual Topics

The screenshot displays the ATU corpus analysis interface, which includes a workflow diagram and several interactive windows.

**Workflow Diagram:** The process starts with a **Corpus** (represented by a document icon), followed by **Preprocess Text** (gear icon), **Topic Modelling** (circular flow icon), **Word Cloud** (cloud icon with a warning triangle), and finally **Corpus Viewer** (magnifying glass icon).

**Topic Modelling Window:** This window shows the results of the topic modelling process. It includes a dropdown for the number of topics (set to 10) and radio buttons for **Latent Semantic Indexing Options** (selected), **Latent Dirichlet Allocation**, and **Hierarchical Dirichlet Process**. The results are displayed in a table:

Topic	Topic keywords
1	said, came, went, little, one, king, go, away, man, father
2	little, hansel, king, gretel, children, mother, bird, red, said, tree
3	bird, hansel, gretel, mother, tree, son, children, beautiful, forest, shudder
4	wife, said, fish, fisherman, came, princess, fox, horse, old, king
5	shudder, youth, wife, fire, fish, learn, fisherman, father, could, boy
6	red, hansel, gretel, rose, white, bear, snow, bird, wolf, cap
7	fox, wolf, pick, bird, ashputtel, bride, tree, cap, home, prince
8	falada, maid, fox, curdken, bride, blow, head, alas, sadly, sparrow
9	cap, grandmother, wolf, bear, white, snow, rose, little, children, hedgehog
10	hedgehog, hare, field, cap, grandmother, run, wife, sparrow, water, little

**Corpus Viewer Window:** This window provides details about the corpus and the selected document. It includes a **RegExp Filter** and a list of documents. The selected document is **A Tale About the Boy Who Went Forth to Learn What Fear Was**. The window also displays a **Word Cloud** for the selected document, showing words like **answered**, **golden**, **might**, **castle**, **began**, **nothing**, **beautiful**, **home**, **set**, **however**, **must**, **ran**, **hansel**, **evening**, **asleep**, **girl**, **live**, **told**, **tell**, **tree**, **like**, **prince**, **thus**, **room**, **great**, **took**, **eat**, **lie**, **old**, **shall**, **well**, **side**, **put**, **red**, **boy**, **went**, **little**, **saw**, **put**, **red**, **boy**, **room**.

LSI provides both positive and negative weights per topic. A positive weight means the word is highly representative of a topic, while a negative weight means the word is highly unrepresentative of a topic (the less it occurs in a text, the more likely the topic). Positive words are colored green and negative words are colored red.

We then select the first topic and display the most frequent words in the topic in **Word Cloud**. We also connected **Preprocess Text** to **Word Cloud** in order to be able to output selected documents. Now we can select a specific word in the word cloud, say *little*. It will be colored red and also highlighted in the word list on the left.

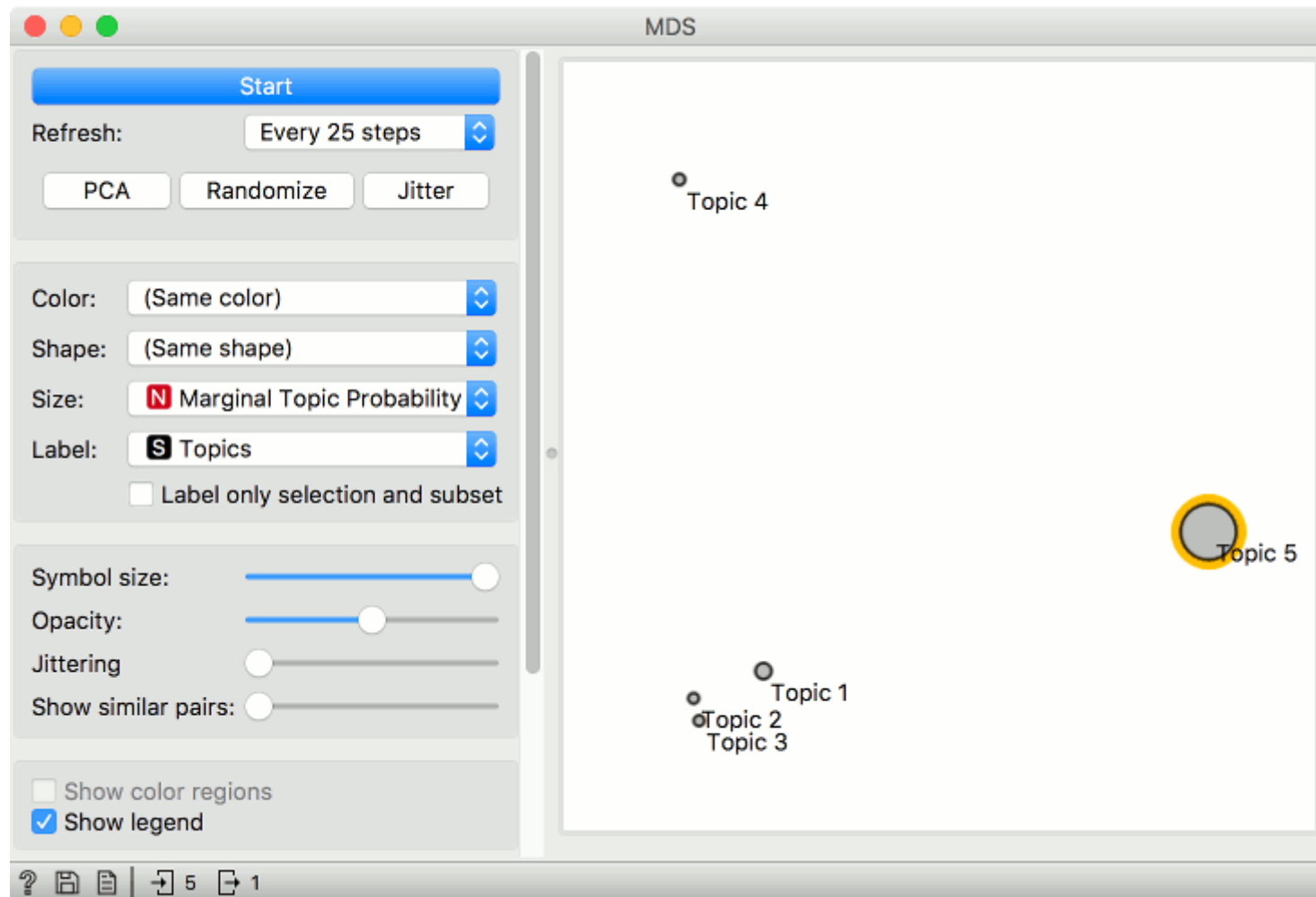
Now we can observe all the documents containing the word *little* in **Corpus Viewer**.

## Topic Visualization

In the second example, we will look at the correlation between topics and words/documents. We are still using the *grimm-tales-selected.tab* corpus. In **Preprocess Text** we are using the default preprocessing, with an additional filter by *document frequency* (0.1 - 0.9). In **Topic Modelling** we are using LDA model with 5 topics.

Connect Topic Modelling to **MDS**. Ensure the link is set to *All Topics - Data*. Topic Modelling will output a matrix of word weights by topic.

In MDS, the points are now topics. We have set the size of the points to *Marginal topic probability*, which is an additional columns of *All Topics* - it reports on the marginal probability of the topic in the corpus (how strongly represented is the topic in the corpus).



We can now explore which words are representative for the topic. Select, say, Topic 5 from the plot and connect MDS to **Box Plot**. Make sure the output is set to *Data - Data* (not *Selected Data - Data*).

In Box Plot, set the subgroup to Selected and check the *Order by relevance to subgroups* box. This option will sort the variables by how well they separate between the selected subgroup values. In our case, this means which words are the most representative for the topic we have selected in the plot (subgroup Yes means selected).

We can see that little, children and kings are the most representative words for Topic 5, with good separation between the word frequency for this topic and all the others. Select other topics in MDS and see how the Box Plot changes.

