
Orange3 Text Mining Documentation

Biolab

Dec 11, 2018

Contents

1	Widgets	1
2	Scripting	45
3	Indices and tables	47

1.1 Corpus

Load a corpus of text documents, (optionally) tagged with categories.

Inputs None

Outputs

Corpus A collection of documents.

Corpus widget reads text corpora from files and sends a corpus instance to its output channel. History of the most recently opened files is maintained in the widget. The widget also includes a directory with sample corpora that come pre-installed with the add-on.

The widget reads data from Excel (**.xlsx**), comma-separated (**.csv**) and native tab-delimited (**.tab**) files.

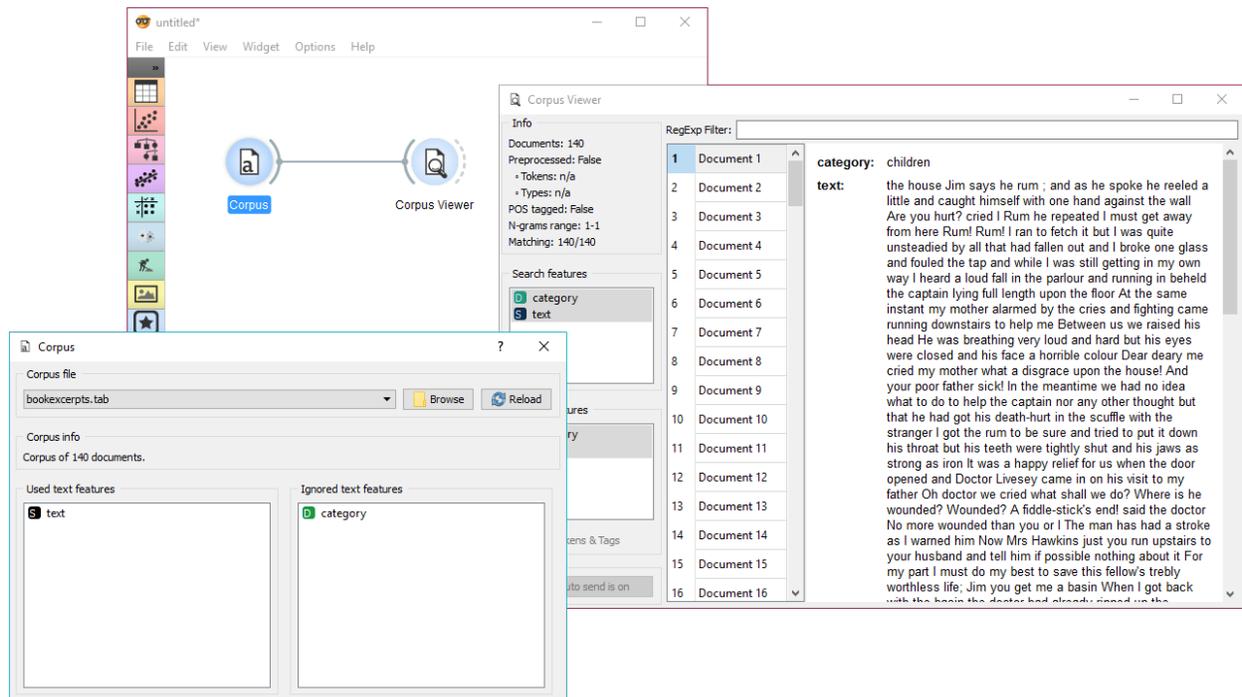
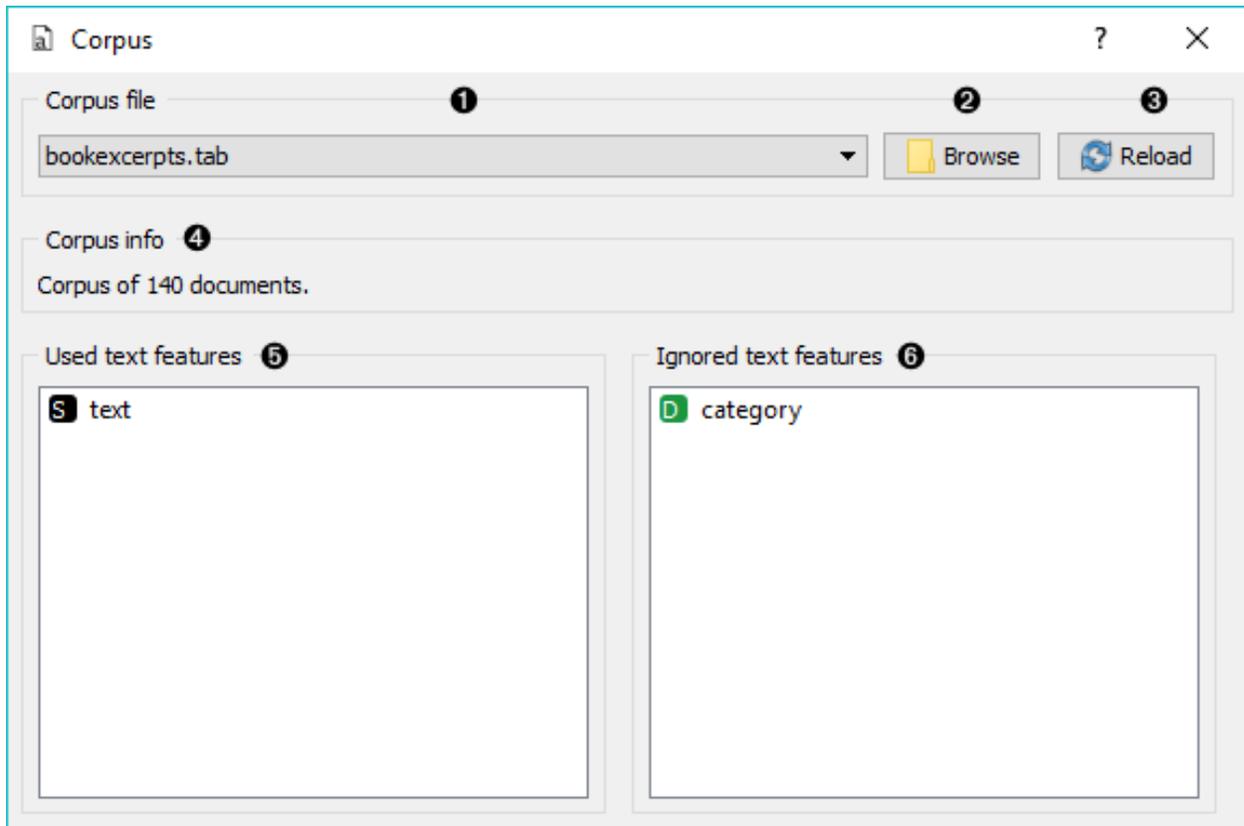
1. Browse through previously opened data files, or load any of the sample ones.
2. Browse for a data file.
3. Reloads currently selected data file.
4. Information on the loaded data set.
5. Features that will be used in text analysis.
6. Features that won't be used in text analysis and serve as labels or class.

You can drag and drop features between the two boxes and also change the order in which they appear.

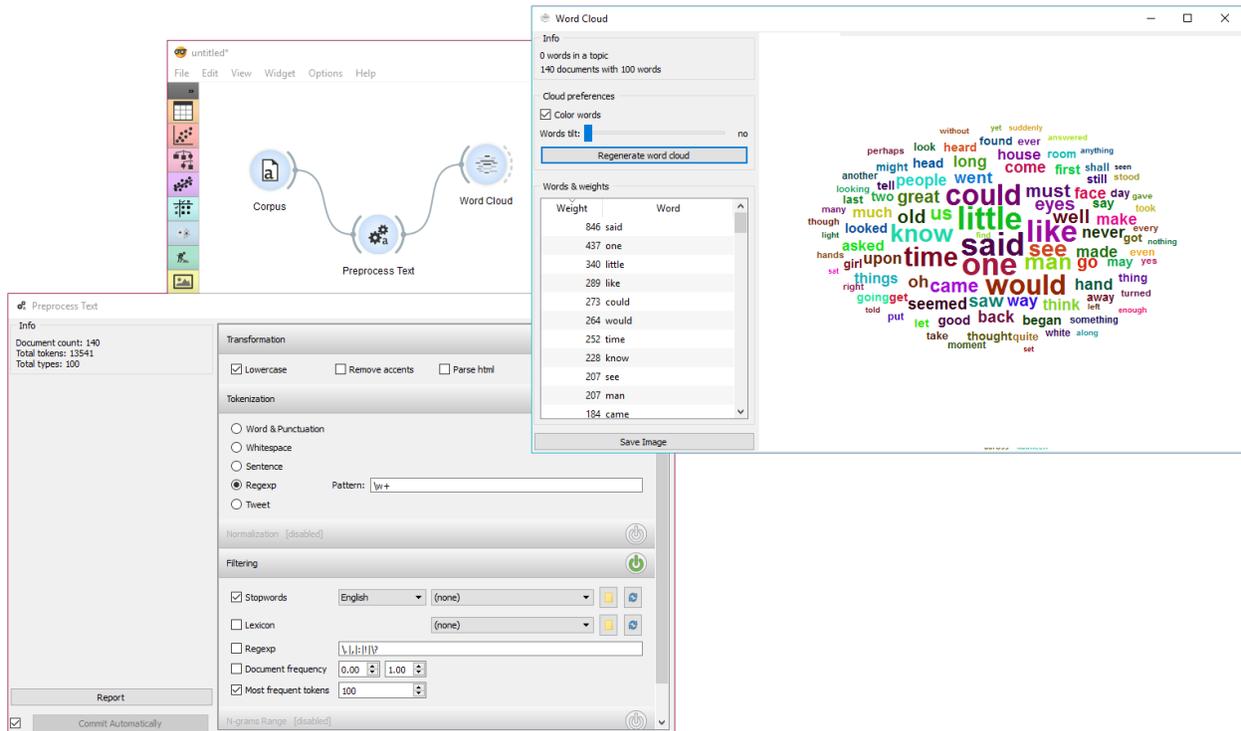
1.1.1 Example

The first example shows a very simple use of **Corpus** widget. Place **Corpus** onto canvas and connect it to *Corpus Viewer*. We've used *booxexcerpts.tab* data set, which comes with the add-on, and inspected it in **Corpus Viewer**.

The second example demonstrates how to quickly visualize your corpus with *Word Cloud*. We could connect **Word Cloud** directly to **Corpus**, but instead we decided to apply some preprocessing with *Preprocess Text*. We are again



working with *book-excerpts.tab*. We've put all text to lowercase, tokenized (split) the text to words only, filtered out English stopwords and selected a 100 most frequent tokens.



1.2 Import Documents

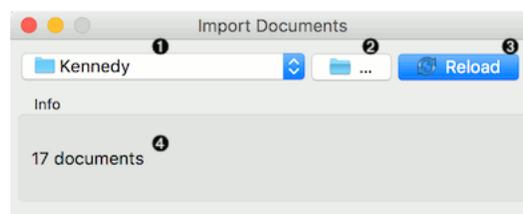
Import text documents from folders.

Inputs None

Outputs

Corpus A collection of documents from the local machine.

Import Documents widget retrieves text files from folders and creates a corpus. The widget reads .txt, .docx, .odt, .pdf and .xml files. If a folder contains subfolders, they will be used as class labels.

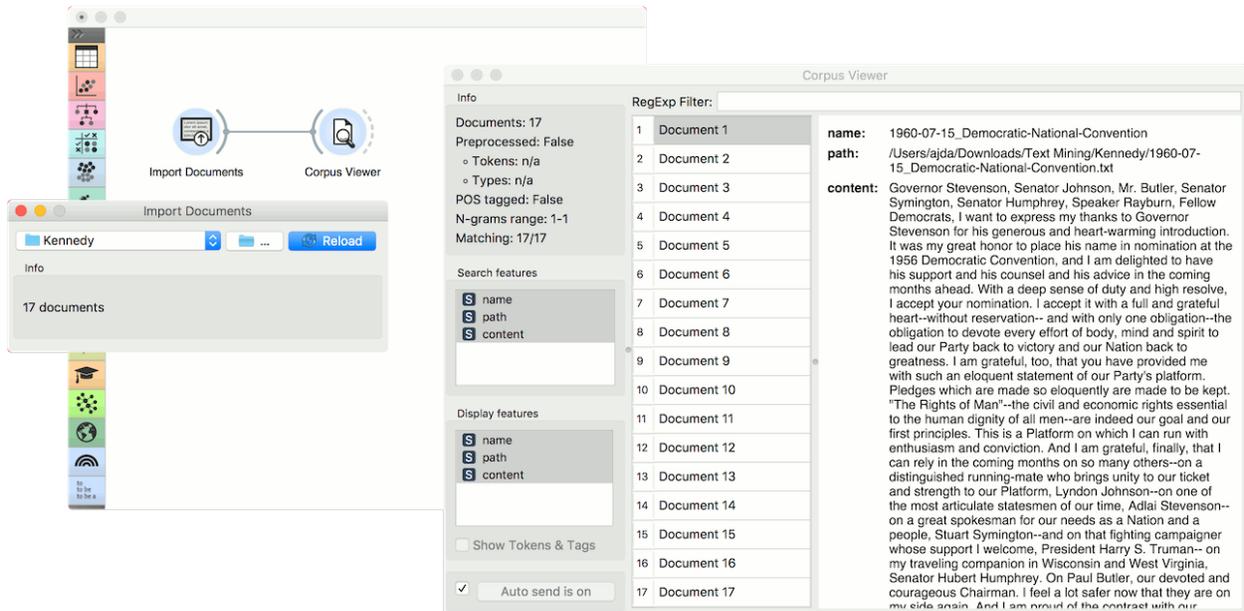


1. Folder being loaded.
2. Load folder from a local machine.
3. Reload the data.
4. Number of documents retrieved.

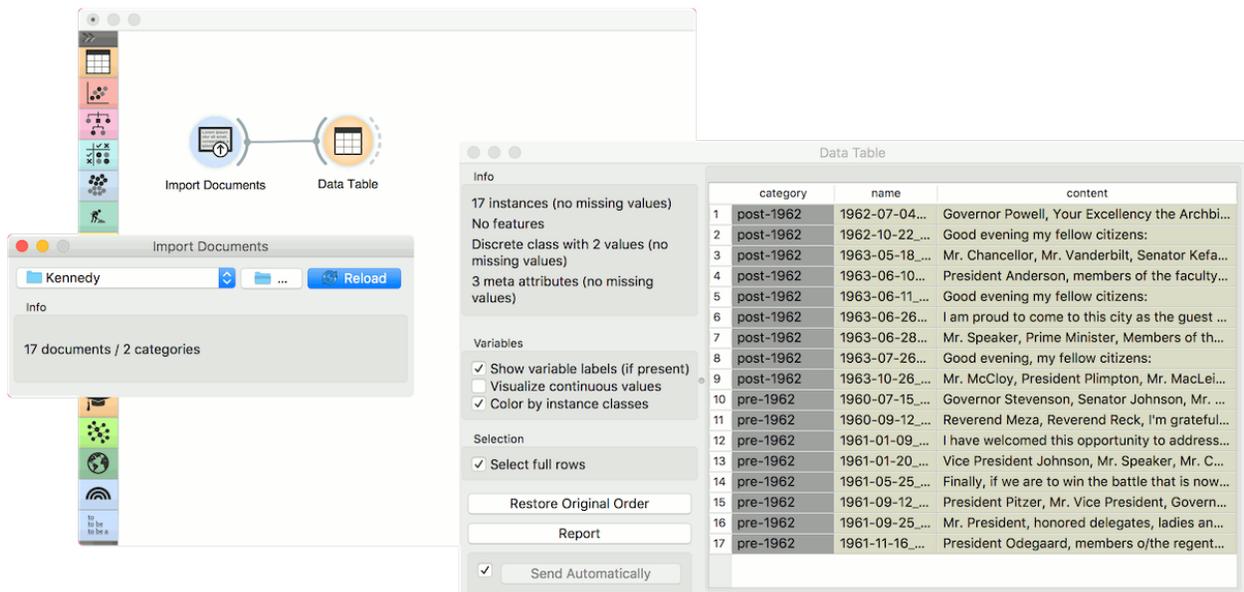
If the widget cannot read the file for some reason, the file will be skipped. Files that were successfully retrieved will still be on the output.

1.2.1 Example

To retrieve the data, select the folder icon on the right side of the widget. Select the folder you wish to turn into corpus. Once the loading is finished, you will see how many documents the widget retrieved. To inspect them, connect the widget to *Corpus Viewer*. We've used a set of Kennedy's speeches in a plain text format.



Now let us try it with subfolders. We have placed Kennedy's speeches in two folders - pre-1962 and post-1962. If I load the parent folder, these two subfolders will be used as class labels. Check the output of the widget in a **Data Table**.



1.3 The Guardian

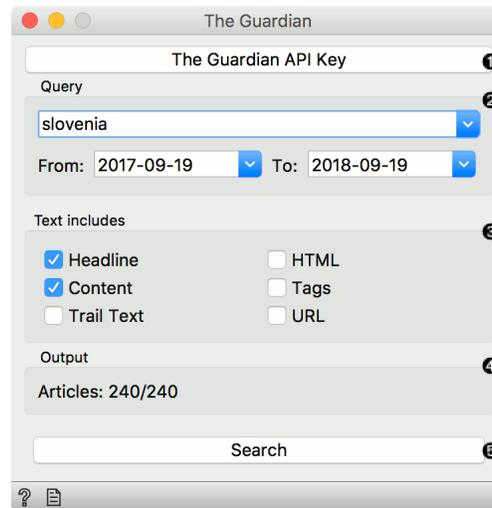
Fetching data from [The Guardian Open Platform](#).

Inputs None

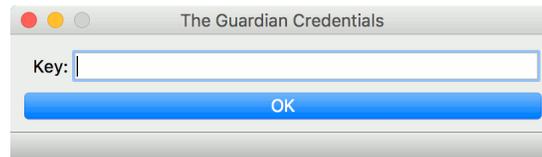
Outputs

Corpus A collection of documents from the Guardian newspaper.

Guardian retrieves articles from the Guardian newspaper via their API. For the widget to work, you need to provide the API key, which you can get at [their access platform](#).



1. Insert the API key for the widget to work.



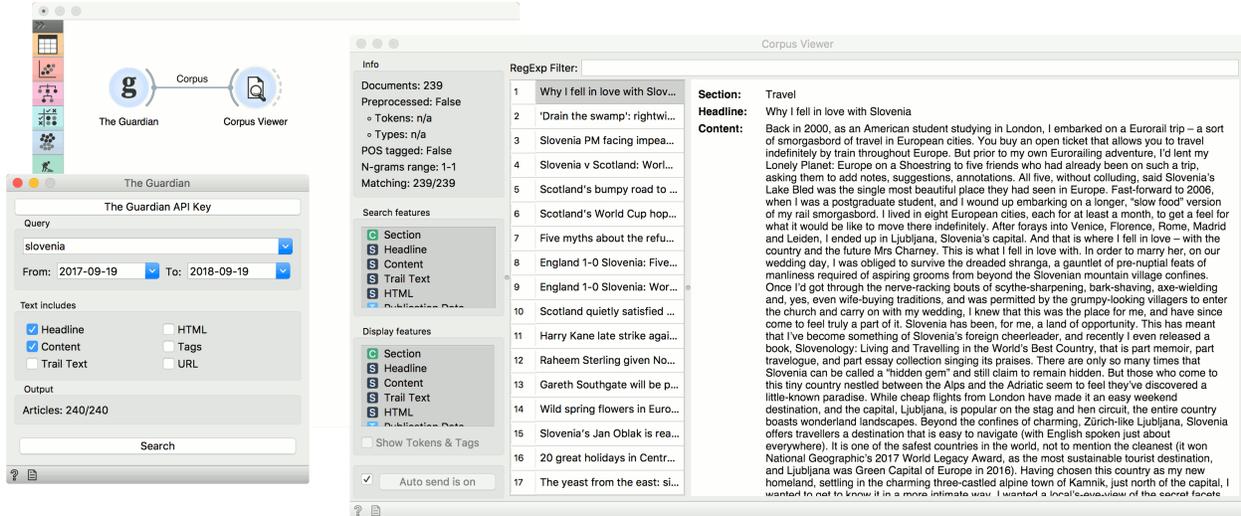
2. Provide the query and set the time frame from which to retrieve the articles.
3. Define which features to retrieve from the Guardian platform.
4. Information on the output.
5. Press *Search* to start retrieving the articles or *Stop* to stop the retrieval.

1.3.1 Example

Guardian can be used just like any other data retrieval widget in Orange, namely *NY Times*, *Wikipedia*, *Twitter* or *PubMed*.

We will retrieve 240 articles mentioning *slovenia* between september 2017 and september 2018. The text will include article headline and content. Upon pressing *Search*, the articles will be retrieved.

We can observe the results in the *Corpus Viewer* widget.



1.4 NY Times

Loads data from the New York Times' [Article Search API](#).

Inputs None

Outputs

Corpus A collection of documents from the New York Times newspaper.

NYTimes widget loads data from New York Times' [Article Search API](#). You can query NYTimes articles from September 18, 1851 to today, but the API limit is set to allow retrieving only a 1000 documents per query. Define which features to use for text mining, *Headline* and *Abstract* being selected by default.

To use the widget, you must enter [your own API key](#).

1. To begin your query, insert NY Times' [Article Search API](#) key. The key is securely saved in your system keyring service (like [Credential Vault](#), [Keychain](#), [KWallet](#), etc.) and won't be deleted when clearing widget settings.
2. **Set query parameters:**
 - *Query*
 - Query time frame. The widget allows querying articles from September 18, 1851 onwards. Default is set to 1 year back from the current date.
3. Define which features to include as text features.
4. Information on the output.
5. Produce report.
6. Run or stop the query.

1.4.1 Example

NYTimes is a data retrieving widget, similar to [Twitter](#) and [Wikipedia](#). As it can retrieve geolocations, that is geographical locations the article mentions, it is great in combination with [GeoMap](#) widget.

First, let's query **NYTimes** for all articles on Slovenia. We can retrieve the articles found and view the results in [Corpus Viewer](#). The widget displays all the retrieved features, but includes on selected features as text mining features.

NY Times (4%, ETA: 0:...

Article API Key ❶

Query ❷

slovenia

From: 2015-10-11 To: 2016-10-10

Text includes ❸

Headline URL

Abstract Locations

Snippet Persons

Lead Paragraph Organizations

Subject Keywords Creative Works

Output ❹

Articles: 20/410

Report ❺ Stop ❻

New York Times API key ? X

Key: |

OK

The screenshot displays the Orange3 Text Mining workflow. On the left, the 'NY Times' widget is configured with a query of 'slovenia' and a date range from 2015-10-11 to 2016-10-10. The 'Text includes' section has 'Headline', 'Abstract', and 'Snippet' checked. The 'Output' section shows 'Articles: 110/410'. The 'Corpus Viewer' widget in the center shows 110 documents with 110/110 matches. A list of search features is visible, including Section, Headline, Abstract, Snippet, Lead Paragraph, Subject Keywords, URL, and Locations. The 'GeoMap' widget on the right shows a map of Europe with Slovenia highlighted in red. A detailed view of a document is shown on the right, including the following text:

Section: World
Headline: Secrecy Reigns as U.N. Seeks a New Secretary General
Abstract: United Nations Security Council straw poll for next secretary general results in former Portuguese Prime Min Antonio Guterres and former Slovenian Pres Danilo Turk receiving most votes, with Bulgarian diplomat Irina Bokova coming in at strong third; proceedings are characterized by secrecy and lack of transparency, for which organization has long been criticized; final choice is ultimately likely to be made by United States or Russia, each of which can veto any candidate.
Snippet: Diplomats emerged from the nearly two-hour voting session with lips sealed about whom they had given a thumb's up or thumb's down to.
Lead Paragraph: Diplomats emerged from the nearly two-hour voting session with lips sealed about whom they had given a thumb's up or thumb's down to.
Subject Keywords: Elections, Diplomatic Service, Embassies and Consulates, United States International Relations
URL: <http://www.nytimes.com/2016/07/22/world/americas/united-nations-secretary-general.html>
Locations: Slovenia, Portugal, Bulgaria, Russia

Now, let's inspect the distribution of geolocations from the articles mentioning Slovenia. We can do this with *GeoMap*. Unsurprisingly, Croatia and Hungary appear the most often in articles on Slovenia (discounting Slovenia itself), with the rest of Europe being mentioned very often as well.

1.5 Pubmed

Fetch data from [PubMed](#) journals.

Inputs None

Outputs

Corpus A collection of documents from the PubMed online service.

PubMed comprises more than 26 million citations for biomedical literature from MEDLINE, life science journals, and online books. The widget allows you to query and retrieve these entries. You can use regular search or construct advanced queries.

1. Enter a valid e-mail to retrieve queries.

2. **Regular search:**

- *Author*: queries entries from a specific author. Leave empty to query by all authors.
- *From*: define the time frame of publication.
- *Query*: enter the query.

Advanced search: enables you to construct complex queries. See [PubMed's website](#) to learn how to construct such queries. You can also copy-paste constructed queries from the website.

3. *Find records* finds available data from PubMed matching the query. Number of records found will be displayed above the button.

4. Define the output. All checked features will be on the output of the widget.

5. Set the number of record you wish to retrieve. Press *Retrieve records* to get results of your query on the output. Below the button is an information on the number of records on the output.

1.5.1 Example

PubMed can be used just like any other data widget. In this example we've queried the database for records on orchids. We retrieved 1000 records and kept only 'abstract' in our meta features to limit the construction of tokens only to this feature.

We used *Preprocess Text* to remove stopword and words shorter than 3 characters (regex `\b\w{1,2}\b`). This will perhaps get rid of some important words denoting chemicals, so we need to be careful with what we filter out. For the sake of quick inspection we only retained longer words, which are displayed by frequency in *Word Cloud*.

1.6 Twitter

Fetching data from [The Twitter Search API](#).

Inputs None

Outputs

Corpus A collection of tweets from the Twitter API.

Pubmed

Email: mail@mail.com ①

Regular search Advanced search ②

Author:

From: 1800-01-01 to: 2016-10-07

Query: orchid

Number of retrievable records for this search query: 1482

Find records ③

Text includes ④

- Authors
- Article title
- Mesh headings
- Abstract
- URL

Retrieve 1000 records from 1482.

Retrieve records ⑤

Number of records retrieved: 1000

The screenshot shows the Orange3 Text Mining workflow. The Pubmed widget is configured with the following settings:

- Regular search: Advanced search
- Email: mail@mail.com
- Author: [empty]
- From: 1800-01-01 to: 2016-10-07
- Query: orchid
- Number of retrievable records for this search query: 1482
- Find records button
- Text includes:
 - Authors
 - Article title
 - Mesh headings
 - Abstract
 - URL
- Retrieve: 1000 records from 1482.
- Retrieve records button
- Number of records retrieved: 1000

The Word Cloud widget displays the following table of words and weights:

Weight	Word
1776	species
1631	orchid
647	orchids
554	two
550	results
544	study
535	plants
492	floral
472	plant
437	genes
436	mycorrhizal
403	populations
397	using
379	fungi
378	gene

The word cloud visualization shows the following prominent words: species, orchid, populations, genes, analysis, pollination, plants, orchids, floral, mycorrhizal, growth, showed, fungi, results, like, studies, pollen, genetic, different, among, identified, evolution, development, pollinators, revealed, sequence, activity, background, flowering, phylogenetic, significant, patterns, provide, suggest, population, different, among, identified, evolution, development, pollinators, revealed, sequence, activity, background, flowering, phylogenetic, significant, patterns, provide, suggest, population, different, among, identified, evolution, development, pollinators, revealed, sequence, activity, background, flowering, phylogenetic, significant, patterns, provide.

Twitter widget enables querying tweets through Twitter API. You can query by content, author or both and accumulate results should you wish to create a larger data set. The widget only supports REST API and allows queries for up to two weeks back.

- To begin your queries, insert Twitter key and secret. They are securely saved in your system keyring service (like Credential Vault, Keychain, KWallet, etc.) and won't be deleted when clearing widget settings. You must first create a [Twitter app](#) to get API keys.
- Set query parameters:**
 - Query word list*: list desired queries, one per line. Queries are automatically joined by OR.
 - Search by*: specify whether you want to search by content, author or both. If searching by author, you must enter proper Twitter handle (without @) in the query list.
 - Allow retweets*: if 'Allow retweets' is checked, retweeted tweets will also appear on the output. This might duplicate some results.
 - Date*: set the query time frame. Twitter only allows retrieving tweets from up to two weeks back.
 - Language*: set the language of retrieved tweets. Any will retrieve tweets in any language.
 - Max tweets*: set the top limit of retrieved tweets. If box is not ticked, no upper bound will be set - widget will retrieve all available tweets.
 - Accumulate results*: if 'Accumulate results' is ticked, widget will append new queries to the previous ones. Enter new queries, run *Search* and new results will be appended to the previous ones.
- Define which features to include as text features.
- Information on the number of tweets on the output.
- Produce report.
- Run query.

Twitter

Twitter API Key 1

Query 2

Query word list:

Multiple lines are automatically joined with OR.

Search by: Content

Allow retweets:

Date: since 2016-09-30 until 2016-10-10

Language: Any

Max tweets: 100

Accumulate results:

Text includes 3

Content Author Description

Info 4

Tweets on output: 0

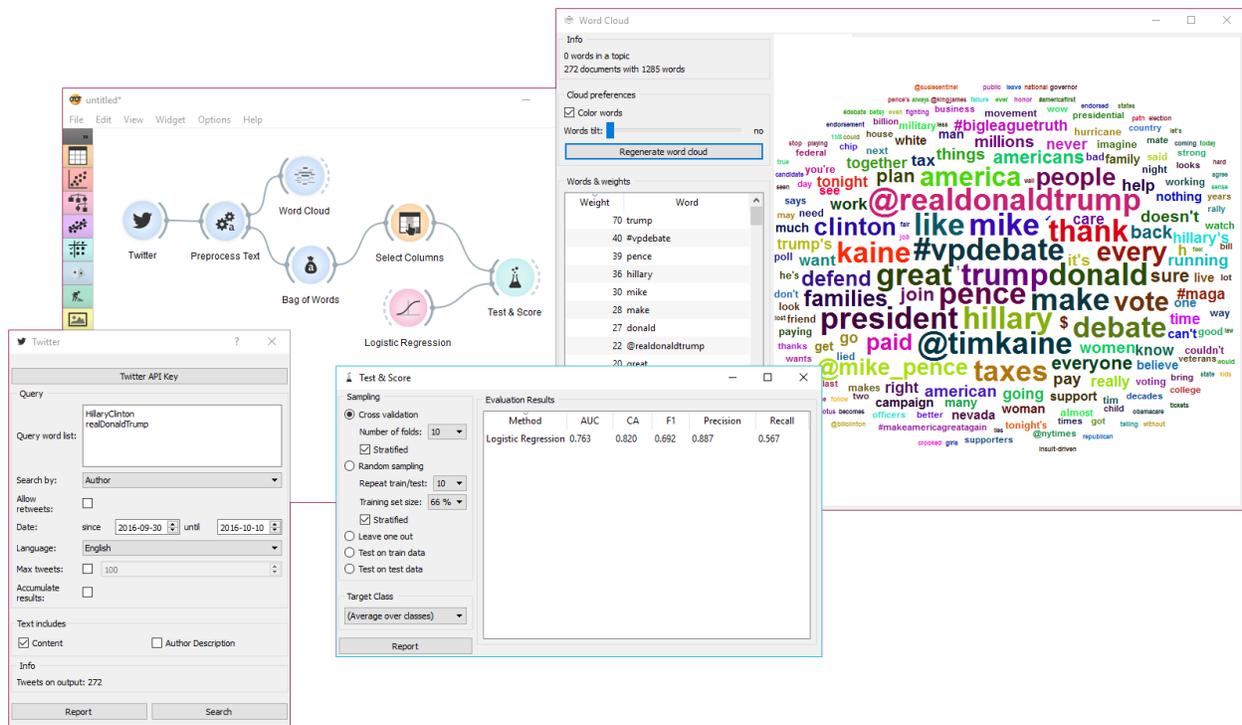
Report 5 Search 6

Twitter API Credentials

Key:

Secret:

OK



1.7 Wikipedia

Fetching data from [MediaWiki RESTful web service API](#).

Inputs None

Outputs

Corpus A collection of documents from the Wikipedia.

Wikipedia widget is used to retrieve texts from Wikipedia API and it is useful mostly for teaching and demonstration.

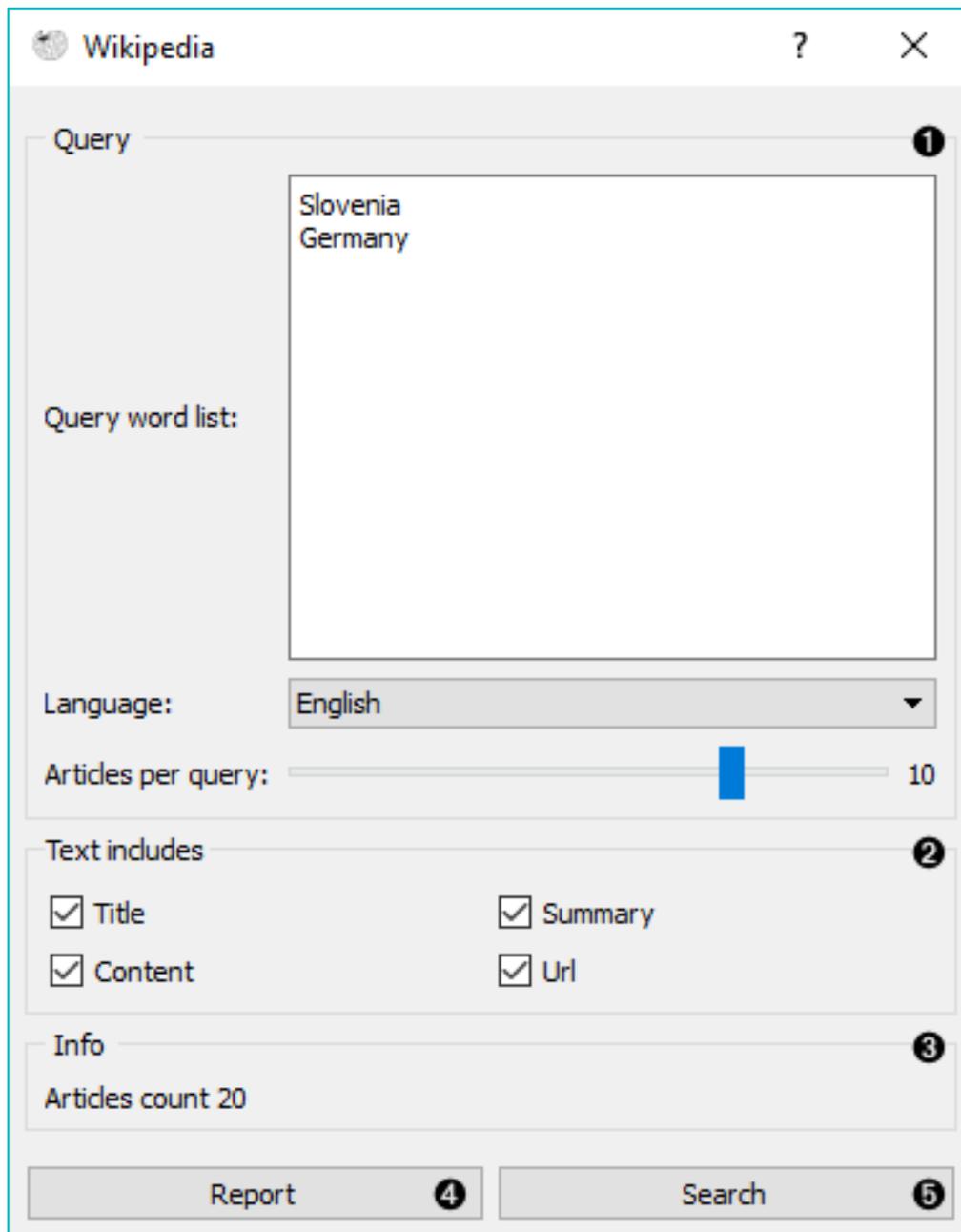
1. Query parameters:

- Query word list, where each query is listed in a new line.
- Language of the query. English is set by default.
- Number of articles to retrieve per query (range 1-25). Please note that querying is done recursively and that disambiguations are also retrieved, sometimes resulting in a larger number of queries than set on the slider.

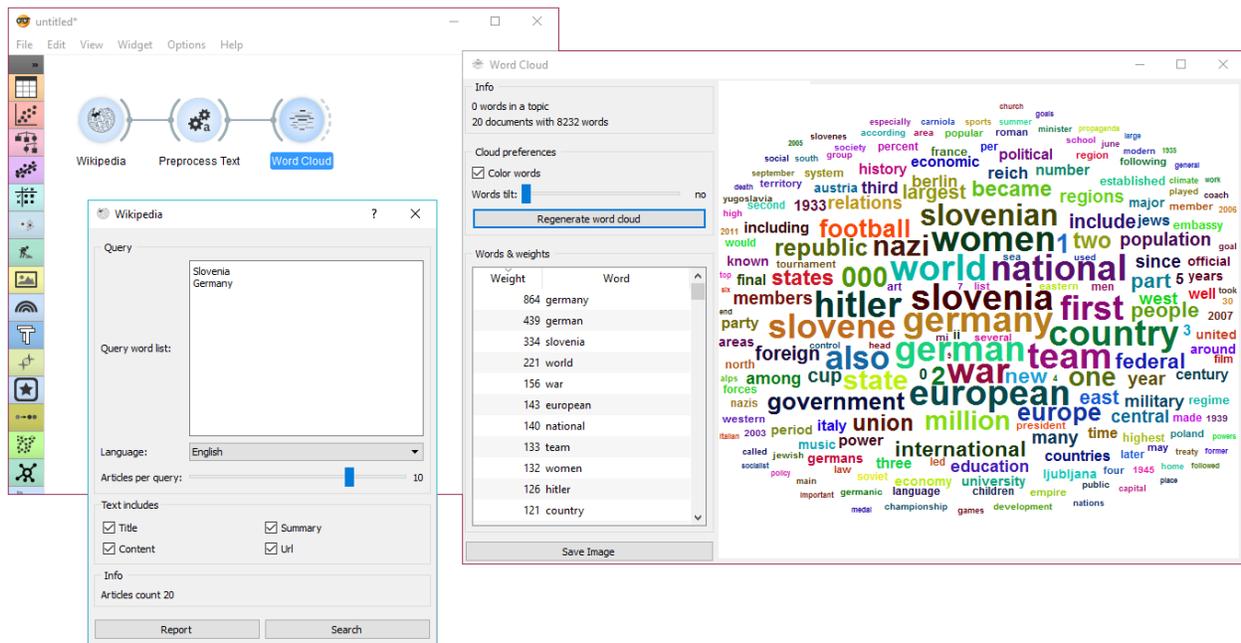
2. Select which features to include as text features.
3. Information on the output.
4. Produce a report.
5. Run query.

1.7.1 Example

This is a simple example, where we use **Wikipedia** and retrieve the articles on ‘Slovenia’ and ‘Germany’. Then we simply apply default preprocessing with *Preprocess Text* and observe the most frequent words in those articles with



Word Cloud.



Wikipedia works just like any other corpus widget (*NY Times*, *Twitter*) and can be used accordingly.

1.8 Preprocess Text

Preprocesses corpus with selected methods.

Inputs

Corpus A collection of documents.

Outputs

Corpus Preprocessed corpus.

Preprocess Text splits your text into smaller units (tokens), filters them, runs *normalization* (stemming, lemmatization), creates *n-grams* and tags tokens with *part-of-speech* labels. Steps in the analysis are applied sequentially and can be turned on or off.

1. **Information on preprocessed data.** *Document count* reports on the number of documents on the input. *Total tokens* counts all the tokens in corpus. *Unique tokens* excludes duplicate tokens and reports only on unique tokens in the corpus.
2. **Transformation transforms input data. It applies lowercase transformation by default.**
 - *Lowercase* will turn all text to lowercase.
 - *Remove accents* will remove all diacritics/accents in text. naïve → naive
 - *Parse html* will detect html tags and parse out text only. <a href...>Some text → Some text
 - *Remove urls* will remove urls from text. This is a <http://orange.biolab.si/> url. → This is a url.
3. **Tokenization is the method of breaking the text into smaller components (words, sentences, bigrams).**
 - *Word & Punctuation* will split the text by words and keep punctuation symbols. This example. → (This), (example), (.)

Preprocess Text

Info

Document count: 140
Total tokens: 64555
Unique tokens: 7392

Transformation

Lowercase Remove accents Parse html Remove urls

Tokenization

Word & Punctuation
 Whitespace
 Sentence
 Regexp Pattern:
 Tweet

Normalization

Porter Stemmer
 Snowball Stemmer Language:
 WordNet Lemmatizer

Filtering

Stopwords
 Lexicon
 Regexp
 Document frequency
 Most frequent tokens

N-grams Range

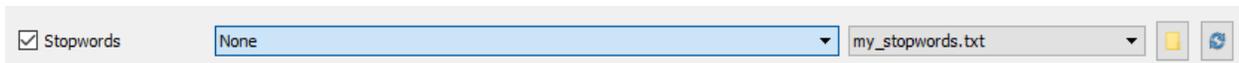
Range:

POS Tagger

Averaged Perceptron Tagger
 Treebank POS Tagger (MaxEnt)
 Stanford POS Tagger

 Commit Automatically

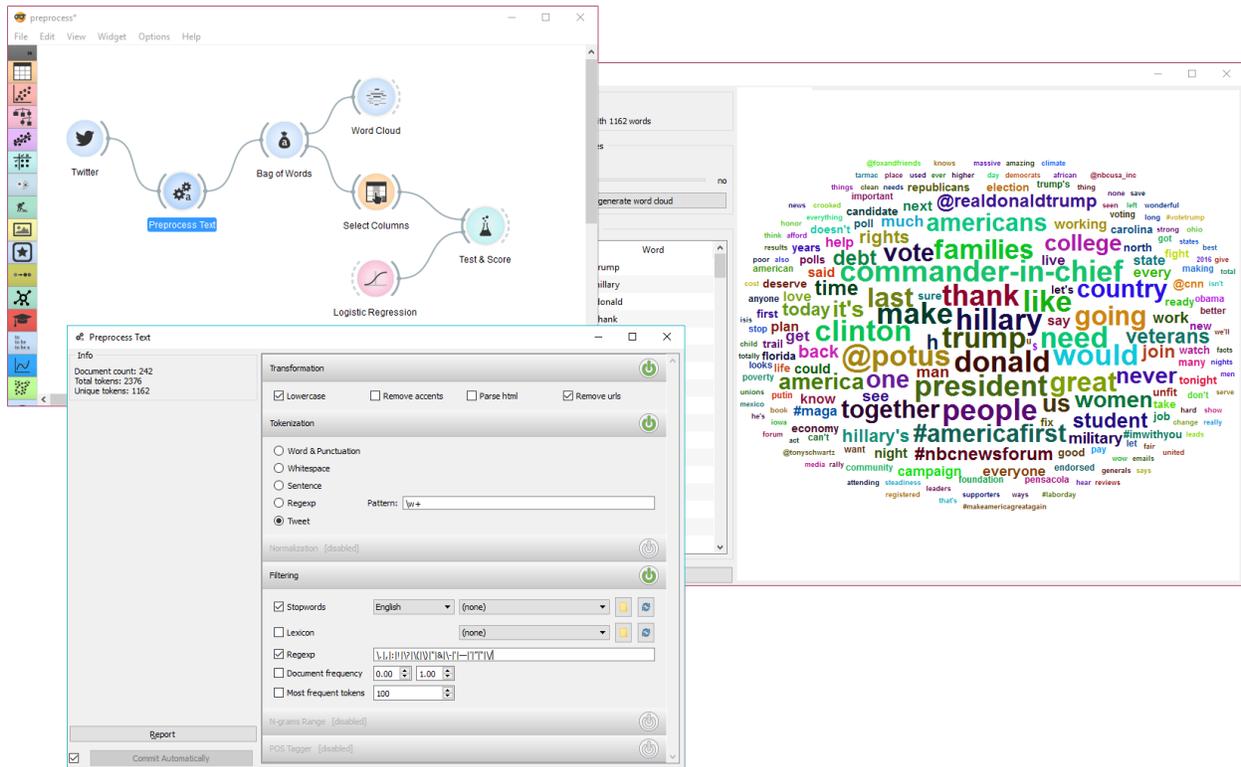
- **Whitespace** will split the text by whitespace only. This example. → (This), (example.)
 - **Sentence** will split the text by fullstop, retaining only full sentences. This example. Another example. → (This example.), (Another example.)
 - **Regex** will split the text by provided regex. It splits by words only by default (omits punctuation).
 - **Tweet** will split the text by pre-trained Twitter model, which keeps hashtags, emoticons and other special symbols. This example. :-) #simple → (This), (example), (.), (:)), (#simple)
4. **Normalization applies stemming and lemmatization to words.** (I've always loved cats. → I have always love cat.) For language:
- **Porter Stemmer** applies the original Porter stemmer.
 - **Snowball Stemmer** applies an improved version of Porter stemmer (Porter2). Set the language for normalization, default is English.
 - **WordNet Lemmatizer** applies a networks of cognitive synonyms to tokens based on a large lexical database of English.
5. **Filtering removes or keeps a selection of words.**
- **Stopwords** removes stopwords from text (e.g. removes 'and', 'or', 'in'...). Select the language to filter by, English is set as default. You can also load your own list of stopwords provided in a simple *.txt file with one stopword per line.



Click 'browse' icon to select the file containing stopwords. If the file was properly loaded, its name will be displayed next to pre-loaded stopwords. Change 'English' to 'None' if you wish to filter out only the provided stopwords. Click 'reload' icon to reload the list of stopwords.

- **Lexicon** keeps only words provided in the file. Load a *.txt file with one word per line to use as lexicon. Click 'reload' icon to reload the lexicon.
 - **Regex** removes words that match the regular expression. Default is set to remove punctuation.
 - **Document frequency** keeps tokens that appear in not less than and not more than the specified number / percentage of documents. If you provide integers as parameters, it keeps only tokens that appear in the specified number of documents. E.g. DF = (3, 5) keeps only tokens that appear in 3 or more and 5 or less documents. If you provide floats as parameters, it keeps only tokens that appear in the specified percentage of documents. E.g. DF = (0.3, 0.5) keeps only tokens that appear in 30% to 50% of documents. Default returns all tokens.
 - **Most frequent tokens** keeps only the specified number of most frequent tokens. Default is a 100 most frequent tokens.
6. **N-grams Range** creates n-grams from tokens. Numbers specify the range of n-grams. Default returns one-grams and two-grams.
7. **POS Tagger runs part-of-speech tagging on tokens.**
- **Averaged Perceptron Tagger** runs POS tagging with Matthew Honnibal's averaged perceptron tagger.
 - **Treebank POS Tagger (MaxEnt)** runs POS tagging with a trained Penn Treebank model.
 - **Stanford POS Tagger** runs a log-linear part-of-speech tagger designed by Toutanova et al. Please download it from the provided website and load it in Orange.
8. Produce a report.

total.



In **Preprocess Text** there's *Tweet* tokenization available, which retains hashtags, emojis, mentions and so on. However, this tokenizer doesn't get rid of punctuation, thus we expanded the Regexp filtering with symbols that we wanted to get rid of. We ended up with word-only tokens, which we displayed in *Word Cloud*. Then we created a schema for predicting author based on tweet content, which is explained in more details in the documentation for *Twitter* widget.

1.9 Bag of Words

Generates a bag of words from the input corpus.

Inputs

Corpus A collection of documents.

Outputs

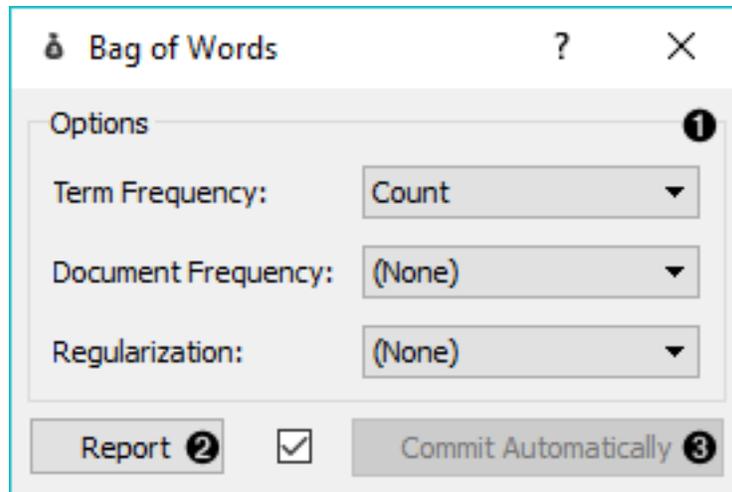
Corpus Corpus with bag of words features appended.

Bag of Words model creates a corpus with word counts for each data instance (document). The count can be either absolute, binary (contains or does not contain) or sublinear (logarithm of the term frequency). Bag of words model is required in combination with *Word Enrichment* and could be used for predictive modelling.

1. Parameters for **bag of words** model:

- **Term Frequency:**

- Count: number of occurrences of a word in a document
- Binary: word appears or does not appear in the document
- Sublinear: logarithm of term frequency (count)



- **Document Frequency:**
 - (None)
 - IDF: inverse document frequency
 - Smooth IDF: adds one to document frequencies to prevent zero division.
- **Regulariation:**
 - (None)
 - L1 (Sum of elements): normalizes vector length to sum of elements
 - L2 (Euclidean): normalizes vector length to sum of squares

2. Produce a report.

3. If *Commit Automatically* is on, changes are communicated automatically. Alternatively press *Commit*.

1.9.1 Example

In the first example we will simply check how the bag of words model looks like. Load *book-excerpts.tab* with *Corpus* widget and connect it to **Bag of Words**. Here we kept the defaults - a simple count of term frequencies. Check what the **Bag of Words** outputs with **Data Table**. The final column in white represents term frequencies for each document.

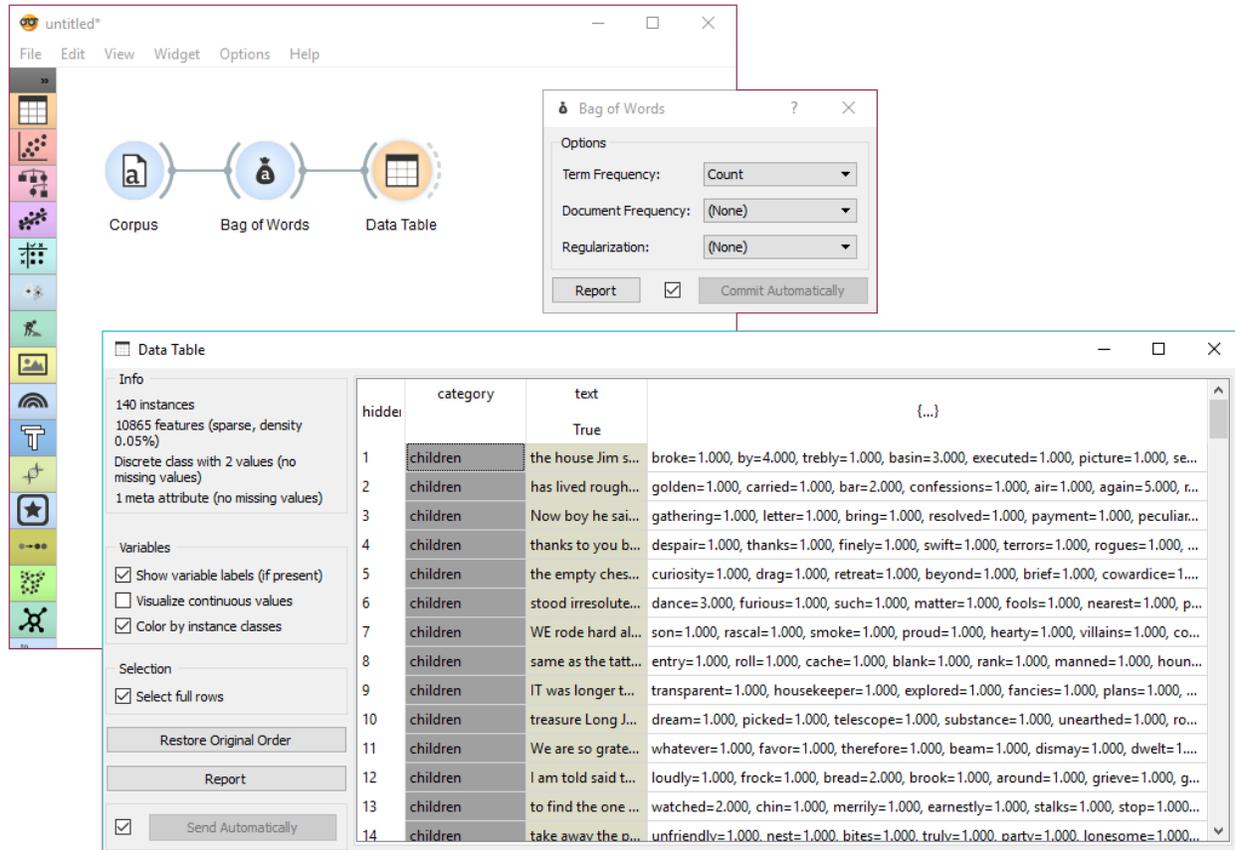
In the second example we will try to predict document category. We are still using the *book-excerpts.tab* data set, which we sent through *Preprocess Text* with default parameters. Then we connected **Preprocess Text** to **Bag of Words** to obtain term frequencies by which we will compute the model.

Connect **Bag of Words** to **Test & Score** for predictive modelling. Connect **SVM** or any other classifier to **Test & Score** as well (both on the left side). **Test & Score** will now compute performance scores for each learner on the input. Here we got quite impressive results with SVM. Now we can check, where the model made a mistake.

Add **Confusion Matrix** to **Test & Score**. Confusion matrix displays correctly and incorrectly classified documents. *Select Misclassified* will output misclassified documents, which we can further inspect with *Corpus Viewer*.

1.10 Similarity Hashing

Computes documents hashes.



Inputs

Corpus A collection of documents.

Outputs

Corpus Corpus with simhash value as attributes.

Similarity Hashing is a widget that transforms documents into similarity vectors. The widget uses **SimHash** from Moses Charikar.

1. Set Simhash size (how many attributes will be on the output, corresponds to bits of information) and shingle length (how many tokens are used in a shingle).
2. *Commit Automatically* output the data automatically. Alternatively, press *Commit*.

1.10.1 Example

We will use *deerwester.tab* to find similar documents in this small corpus. Load the data with *Corpus* and pass it to **Similarity Hashing**. We will keep the default hash size and shingle length. We can observe what the widget outputs in a **Data Table**. There are 64 new attributes available, corresponding to the *Simhash size* parameter.

1.10.2 References

Charikar, M. (2002) Similarity estimation techniques from rounding algorithms. STOC '02 Proceedings of the thirty-fourth annual ACM symposium on Theory of computing, p. 380-388.

The screenshot displays a workflow in Orange3. The main workflow consists of the following widgets: Corpus, Preprocess Text, Bag of Words, SVM, Test & Score, Confusion Matrix, and Corpus Viewer.

Bag of Words Widget Options:

- Term Frequency: Count
- Document Frequency: IDF
- Regularization: (None)

Test & Score Widget Evaluation Results:

Method	AUC	CA	F1	Precision	Recall
SVM	0.971	0.971	0.971	1.000	0.943

Confusion Matrix Widget Output:

		Predicted		Σ
		adult	children	
Actual	adult	70	0	70
	children	4	66	70
Σ		74	66	140

Corpus Viewer Widget Output:

Documents: 4
Preprocessed: False
Tokens: n/a
Types: n/a
POS tagged: False
Ngrams range: 1-1
Matchings: 4/4

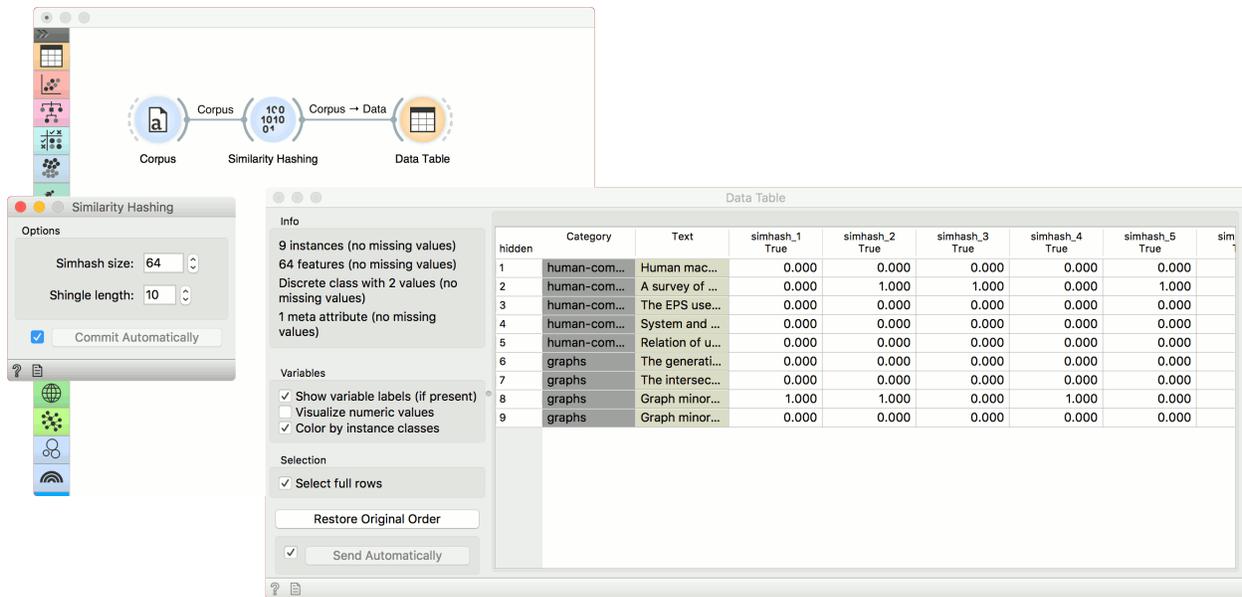
Search features: category, text, category(SVM)
Display features: category, text, category(SVM)

RegExp Filter:

category: children
text: thanks to you big hulking chicken-hearted men We'll have that chest open if we die for it And I'll thank you for that bag Mrs Crossley to bring back our lawful money in Of course I said I would go with my mother and of course they all cried out at our foolishness but even then not a man would go along with us All they would do was to give me a loaded pistol lest we were attacked and to promise to have horses ready saddled in case we were pursued on our return while one lad was to ride forward to the doctor's in search of armed assistance My heart was beating finely when we two set forth in the cold night upon this dangerous venture A full moon was beginning to rise and peered redly through the upper edges of the fog and this increased our haste for it was plain before we came forth again that all would be as bright as day and our departure exposed to the eyes of any watchers We slipped along the hedges noiseless and swift nor did we see or hear anything to increase our terrors till to our relief the door of the Admiral Benbow had closed behind us I slipped the bolt at once and we stood and panted for a moment in the dark alone in the house with the dead captain's body Then my mother got a candle in the bar and holding each

Similarity Hashing Widget Options:

- Simhash size: 64
- Shingle length: 10
- Commit Automatically



1.11 Sentiment Analysis

Predict sentiment from text.

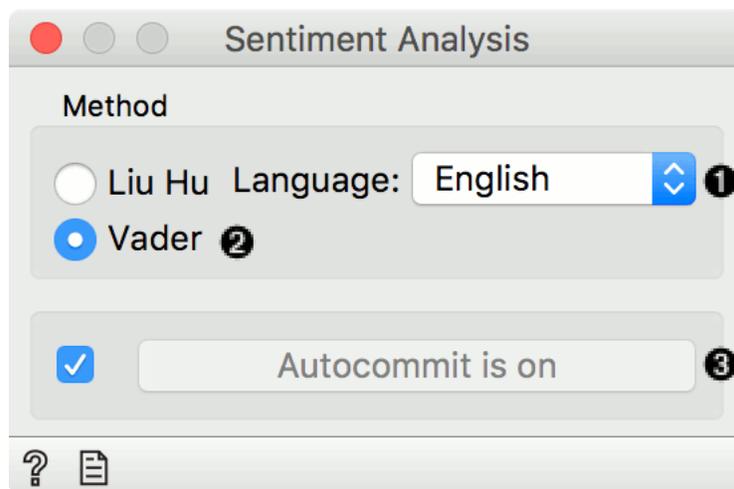
Inputs

Corpus A collection of documents.

Outputs

Corpus A corpus with information on the sentiment of each document.

Sentiment Analysis predicts sentiment for each document in a corpus. It uses Liu Hu and Vader sentiment modules from [NLTK](#). Both of them are lexicon-based. For Liu Hu, you can choose English or Slovenian version.



1. *Method:*

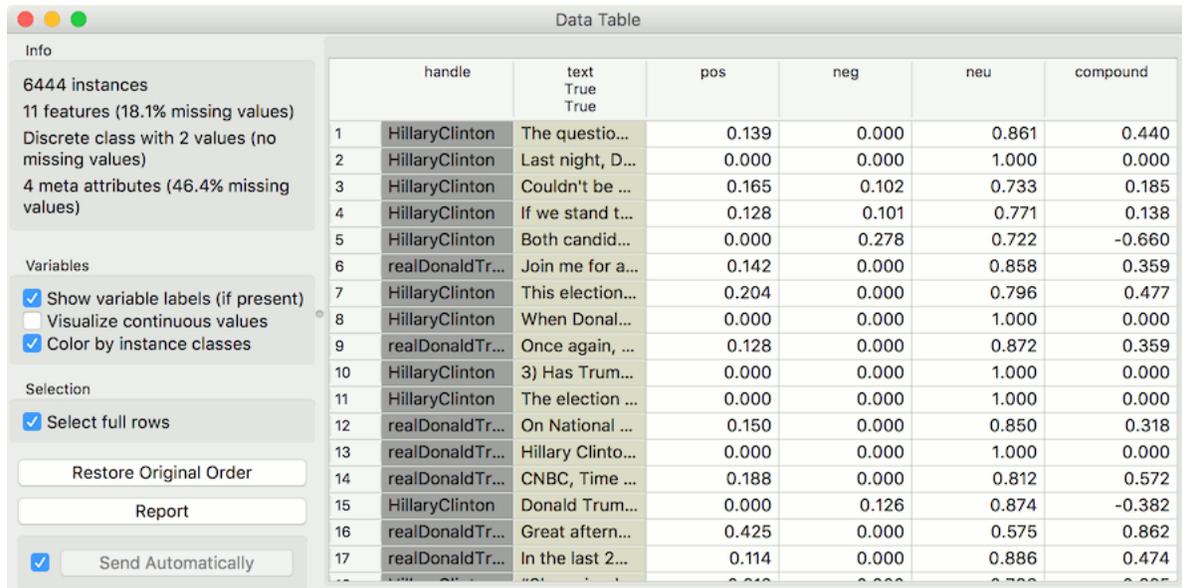
- *Liu Hu*: lexicon-based sentiment analysis (supports English and Slovenian)
- *Vader*: lexicon- and rule-based sentiment analysis

2. Produce a report.
3. If *Auto commit is on*, sentiment-tagged corpus is communicated automatically. Alternatively press *Commit*.

1.11.1 Example

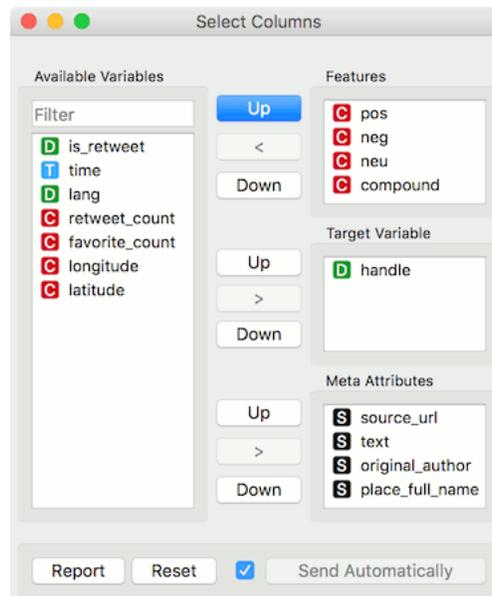
Sentiment Analysis can be used for constructing additional features with sentiment prediction from corpus. First, we load *Election-2016-tweets.tab* in *Corpus*. Then we connect **Corpus** to **Sentiment Analysis**. The widget will append 4 new features for Vader method: positive score, negative score, neutral score and compound (combined score).

We can observe new features in a **Data Table**, where we sorted the *compound* by score. Compound represents the total sentiment of a tweet, where -1 is the most negative and 1 the most positive.

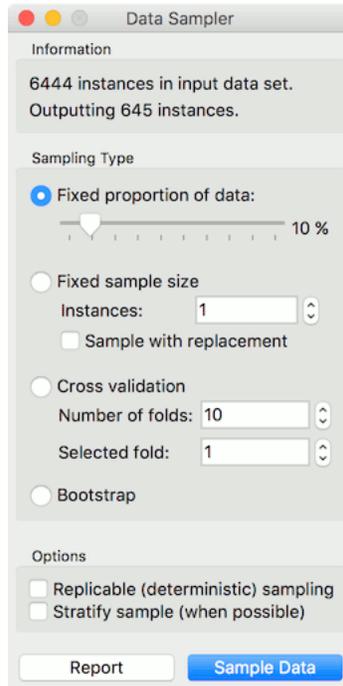


	handle	text True True	pos	neg	neu	compound
1	HillaryClinton	The questio...	0.139	0.000	0.861	0.440
2	HillaryClinton	Last night, D...	0.000	0.000	1.000	0.000
3	HillaryClinton	Couldn't be ...	0.165	0.102	0.733	0.185
4	HillaryClinton	If we stand t...	0.128	0.101	0.771	0.138
5	HillaryClinton	Both candid...	0.000	0.278	0.722	-0.660
6	realDonaldTrump	Join me for a...	0.142	0.000	0.858	0.359
7	HillaryClinton	This election...	0.204	0.000	0.796	0.477
8	HillaryClinton	When Donal...	0.000	0.000	1.000	0.000
9	realDonaldTrump	Once again, ...	0.128	0.000	0.872	0.359
10	HillaryClinton	3) Has Trum...	0.000	0.000	1.000	0.000
11	HillaryClinton	The election ...	0.000	0.000	1.000	0.000
12	realDonaldTrump	On National ...	0.150	0.000	0.850	0.318
13	realDonaldTrump	Hillary Clinto...	0.000	0.000	1.000	0.000
14	realDonaldTrump	CNBC, Time ...	0.188	0.000	0.812	0.572
15	HillaryClinton	Donald Trum...	0.000	0.126	0.874	-0.382
16	realDonaldTrump	Great aftern...	0.425	0.000	0.575	0.862
17	realDonaldTrump	In the last 2...	0.114	0.000	0.886	0.474

Now let us visualize the data. We have some features we are currently not interested in, so we will remove them with **Select Columns**.



Then we will make our corpus a little smaller, so it will be easier to visualize. Pass the data to **Data Sampler** and retain a random 10% of the tweets.



Now pass the filtered corpus to **Heat Map**. Use *Merge by k-means* to merge tweets with the same polarity into one line. Then use *Cluster by rows* to create a clustered visualization where similar tweets are grouped together. Click on a cluster to select a group of tweets - we selected the negative cluster.

To observe the selected subset, pass the tweets to *Corpus Viewer*.

1.11.2 References

Hutto, C.J. and E. E. Gilbert (2014). VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. Eighth International Conference on Weblogs and Social Media (ICWSM-14). Ann Arbor, MI, June 2014.

Hu, Mingqing and Bing Liu (2004). Mining opinion features in customer reviews. In Proceedings of AAAI Conference on Artificial Intelligence, vol. 4, pp. 755–760. [Available online](#).

Kadunc, Klemen and Marko Robnik-Šikonja (2016). Analiza mnenj s pomočjo strojnega učenja in slovenskega leksikona sentimenta. Conference on Language Technologies & Digital Humanities, Ljubljana (in Slovene). [Available online](#).

1.12 Tweet Profiler

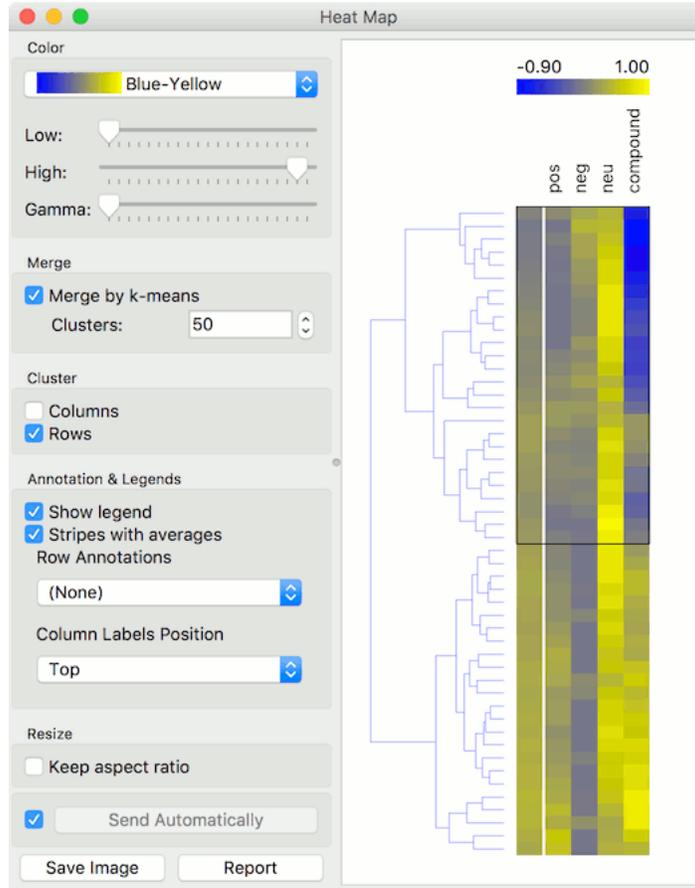
Detect Ekman's, Plutchik's or Profile of Mood States' emotions in tweets.

Inputs

Corpus A collection of tweets (or other documents).

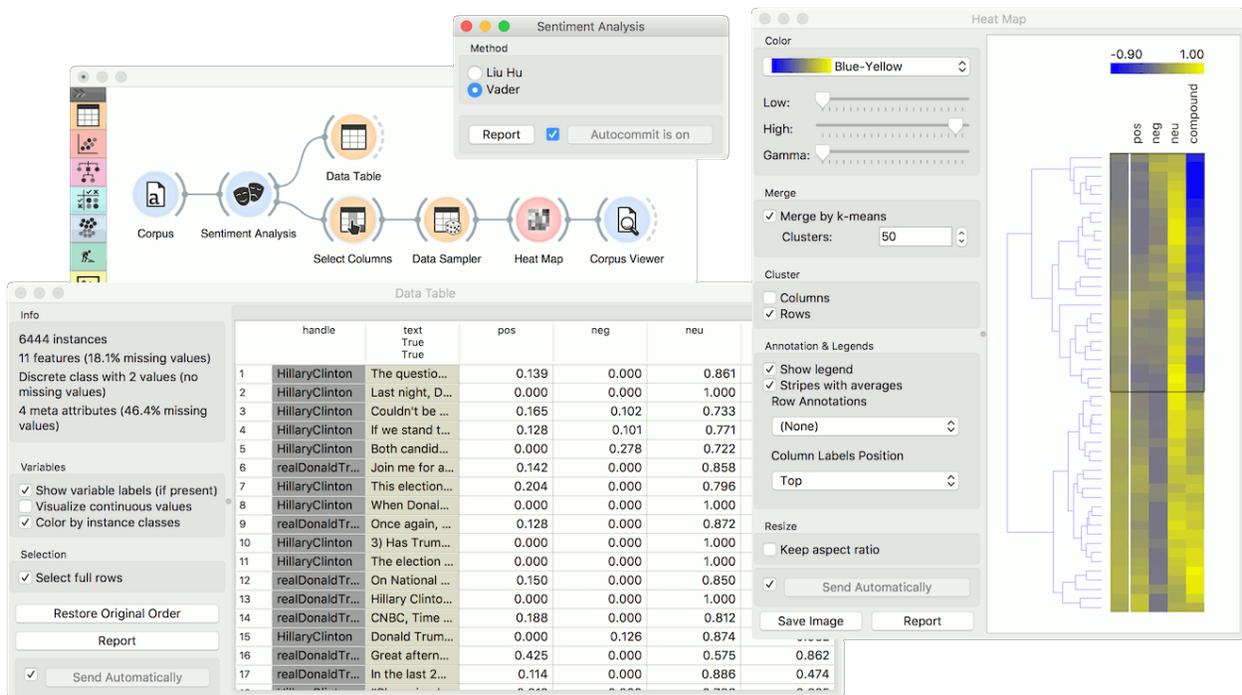
Outputs

Corpus A corpus with information on the sentiment of each document.

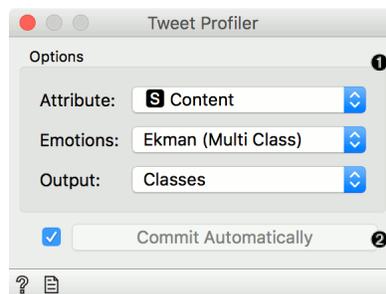


The Corpus Viewer window shows a list of 17 tweets on the left, each with a sentiment score. The right pane provides a detailed view of the selected tweet, including its sentiment score and the full text. The interface includes an 'Info' section on the left with document statistics (375 documents, 375/375 matching), search features (pos, neg, neu, compound, handle), and display features (handle, source_url, text, original_author, place_full_name). A 'RegExp Filter' field is at the top, and an 'Auto send is on' checkbox is at the bottom left.

Index	Tweet Text	Sentiment Score
1	A message of ...	-0.710
2	Let's ask ours...	-0.572
3	Wonderful @p...	-0.611
4	Our diversity i...	-0.599
5	Little Marco R...	-0.887
6	Your @GOP pr...	-0.932
7	.@realDonaldTrump...	-0.832
8	You're wrong, ...	
9	This week, Tru...	
10	Hillary Clinton ...	
11	Leaked e-mail...	
12	The NRA is ba...	
13	Kasich voted f...	
14	Gun violence a...	
15	Wow, the ridic...	
16	So many in the...	
17	If @TedCruz d...	



Tweet Profiler retrieves information on sentiment from the server for each given tweet (or document). The widget sends data to the server, where a model computes emotion probabilities and/or scores. The widget support three classifications of emotion, namely Ekman's, Plutchik's and Profile of Mood States (POMS).



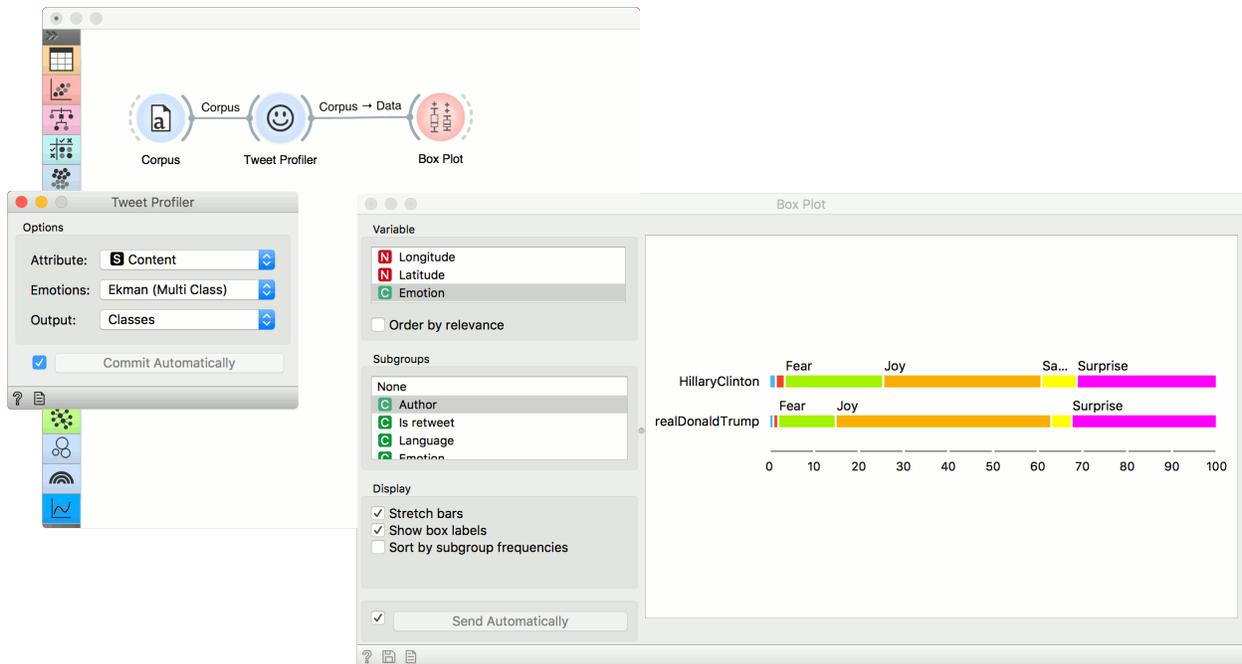
1. Options:

- Attribute to use as content.
- Emotion classification, either Ekman's, Plutchik's or Profile of Mood States. Multi-class will output one most probable emotion per document, while multi-label will output values in columns per each emotion.
- The widget can output classes of emotion (categorical), probabilities (numeric), or embeddings (an emotional vector of the document).

2. *Commit Automatically* automatically outputs the result. Alternatively, press *Commit*.

1.12.1 Example

We will use *election-tweets-2016.tab* for this example. Load the data with *Corpus* and connect it to **Tweet Profiler**. We will use *Content* attribute for the analysis, Ekman's classification of emotion with multi-class option and we will output the result as class. We will observe the results in a **Box Plot**. In the widget, we have selected to observe the *Emotion* variable, grouped by *Author*. This way we can see which emotion prevails by which author.



1.12.2 References

Colnerič, Niko and Janez Demšar (2018). Emotion Recognition on Twitter: Comparative Study and Training a Unison Model. In IEEE Transactions on Affective Computing. [Available online](#).

1.13 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.

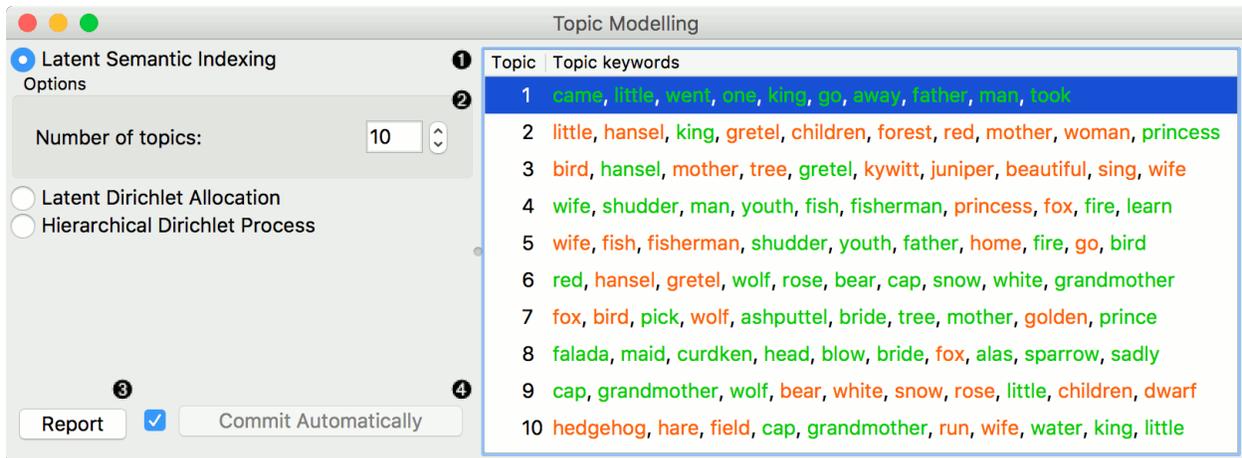
All Topics Topic weights by tokens.

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.

1. Topic modelling algorithm:

- Latent Semantic Indexing
- 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, h



- First level concentration (γ): distribution at the first (corpus) level of Dirichlet Process
- Second level concentration (α): distribution at the second (document) level of Dirichlet Process
- The topic Dirichlet (α): 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 (κ): step size
- Slow down parameter (τ)

3. Produce a report.

4. If *Commit Automatically* is on, changes are communicated automatically. Alternatively press *Commit*.

1.13.1 Example

In the first example, we present a simple use of the **Topic Modelling** widget. First we load *grimm-tales-selected.tab* data set and use *Preprocess Text* to tokenize by words only and remove stopwords. Then we connect **Preprocess Text** to **Topic Modelling**, where we use a simple *Latent Semantic Indexing* to find 10 topics in the text.

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*.

In the second example, we will look at the correlation between topics and words/documents. Connect **Topic Modelling** to **Heat Map**. Ensure the link is set to *All Topics - Data*. **Topic Modelling** will output a matrix of topic weights by words from text (more precisely, tokens).

We can observe the output in a **Data Table**. Tokens are in rows and retrieved topics in columns. Values represent how much a word is represented in a topic.

To visualize this matrix, open **Heat Map**. Select *Merge by k-means* and *Cluster - Rows* to merge similar rows into one and sort them by similarity, which makes the visualization more compact.

The screenshot displays a workflow in Orange3. The 'Topic Modelling' widget is the central focus, showing 10 topics with their respective keywords. A 'Word Cloud' widget is connected to it, displaying a visualization of words from the corpus. The 'Corpus Viewer' widget on the right shows a list of documents and their metadata, including titles and abstracts.

Topic Modelling - Topic keywords:

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

Word Cloud - Top Words:

Weight	Word
0.514	said
0.213	came
0.210	went
0.208	little
0.164	one
0.154	king
0.144	go
0.126	away

Corpus Viewer - Document List:

Doc ID	Title
1	A Tale About the Boy Who We...
2	Brier Rose
3	Cat and Mouse in Partnership
4	Cinderella
5	Hansel and Gretel
6	Little Red Riding Hood
7	Mother Holle
8	Old Sultan
9	Pack of Scoundrels

Data Table

Info: 3716 instances (no missing values), 10 features (no missing values), No target variable, 1 meta attribute (no missing values).

Variables: Show variable labels (if present), Visualize continuous values, Color by instance classes.

Selection: Select full rows.

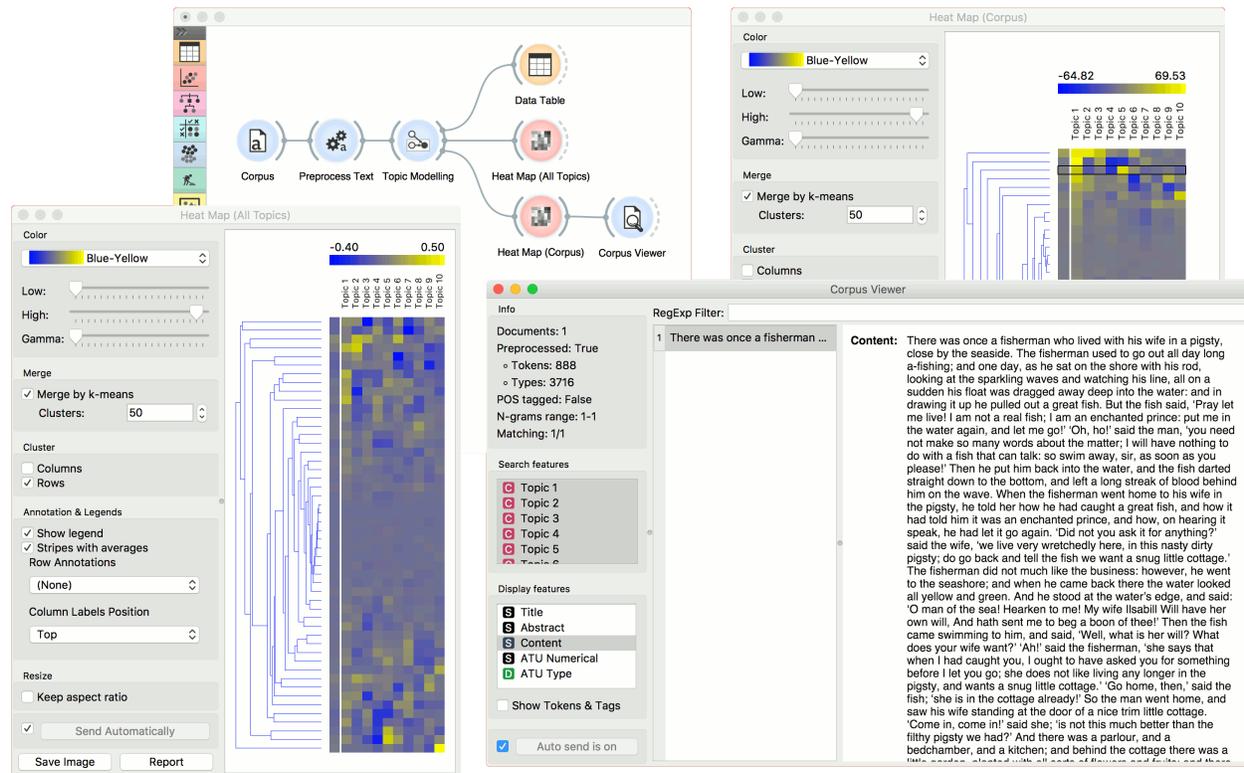
Buttons: Restore Original Order, Report, Send Automatically.

Word	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8	Topic 9	Topic 10
1 _jug	0.000	-0.000	0.000	-0.000	0.000	-0.000	0.000	0.001	0.000	0.000
2 _my_	0.001	-0.001	-0.003	0.003	0.002	0.001	0.010	-0.006	-0.005	0.003
3 abide	0.002	-0.003	0.001	0.004	0.001	0.001	0.000	-0.003	-0.001	-0.003
4 able	0.017	-0.001	0.011	-0.001	-0.015	-0.004	-0.024	-0.009	-0.014	-0.011
5 aboard	0.000	0.000	0.000	-0.000	0.000	-0.000	-0.000	0.000	-0.001	-0.000
6 abode	0.002	-0.003	0.001	0.003	0.001	0.001	0.000	-0.002	-0.001	-0.003
7 abominably	0.000	-0.000	0.000	0.001	-0.000	-0.001	-0.003	-0.000	-0.001	0.001
8 absence	0.000	-0.000	0.000	0.000	0.000	-0.000	-0.000	0.002	-0.000	0.000
9 abundance	0.000	-0.000	0.000	0.001	-0.000	-0.001	-0.003	-0.000	-0.001	0.001
10 accept	0.000	0.000	0.000	0.000	-0.000	-0.000	0.000	0.000	-0.000	-0.000
11 accepted	0.000	0.000	-0.000	-0.002	0.002	0.000	-0.002	0.001	-0.004	0.016
12 accomplish	0.001	-0.001	0.001	-0.002	0.001	-0.001	-0.005	-0.000	-0.004	0.016
13 accomplished	0.000	-0.000	0.000	-0.001	-0.001	0.000	0.000	-0.001	0.000	0.000
14 accord	0.003	-0.003	0.005	-0.008	-0.008	-0.002	0.000	-0.001	-0.001	0.000
15 according	0.001	0.000	0.001	0.001	0.001	-0.000	0.000	0.002	-0.000	0.000
16 accordingly	0.001	-0.001	0.002	-0.000	-0.001	-0.002	-0.004	0.000	-0.001	-0.000
17 account	0.004	-0.002	0.000	0.005	-0.001	-0.009	-0.009	-0.010	0.013	0.003
18 accused	0.001	0.002	0.000	0.001	0.003	-0.010	0.002	-0.003	0.010	0.002
19 accursed	0.001	-0.001	-0.000	0.000	-0.000	0.000	0.001	0.000	-0.000	-0.001
20 accustomed	0.000	0.000	-0.000	-0.002	0.002	0.000	-0.002	0.001	-0.004	0.016

In the upper part of the visualization, we have words that highly define topics 1-3 and in the lower part those that define topics 5 and 10.

We can similarly observe topic representation across documents. We connect another **Heat Map** to **Topic Modelling** and set link to *Corpus - Data*. We set *Merge* and *Cluster* as above.

In this visualization we see how much is a topic represented in a document. Looks like Topic 1 is represented almost across the entire corpus, while other topics are more specific. To observe a specific set of document, select either a clustering node or a row in the visualization. Then pass the data to *Corpus Viewer*.



1.14 Corpus Viewer

Displays corpus content.

Inputs

Corpus A collection of documents.

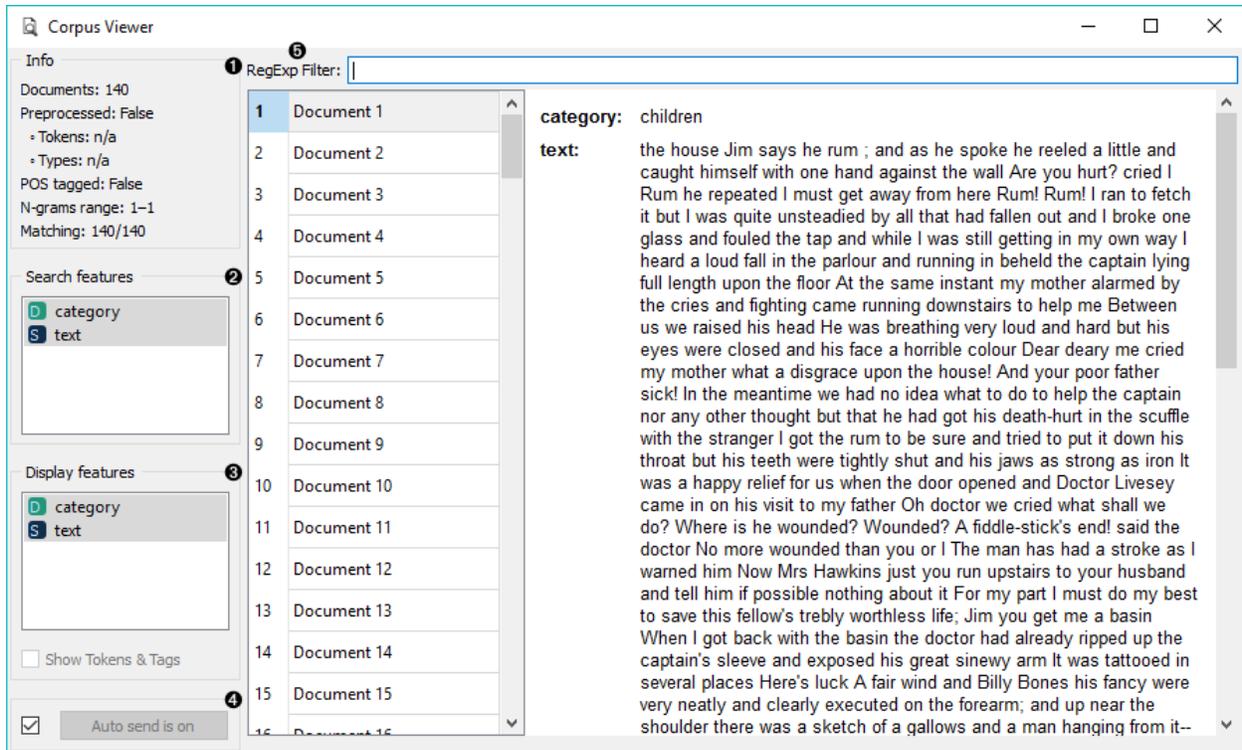
Outputs

Corpus Documents containing the queried word.

Corpus Viewer is meant for viewing text files (instances of *Corpus*). It will always output an instance of corpus. If *RegExp* filtering is used, the widget will output only matching documents.

1. Information:

- *Documents*: number of documents on the input
- *Preprocessed*: if preprocessor is used, the result is True, else False. Reports also on the number of tokens and types (unique tokens).

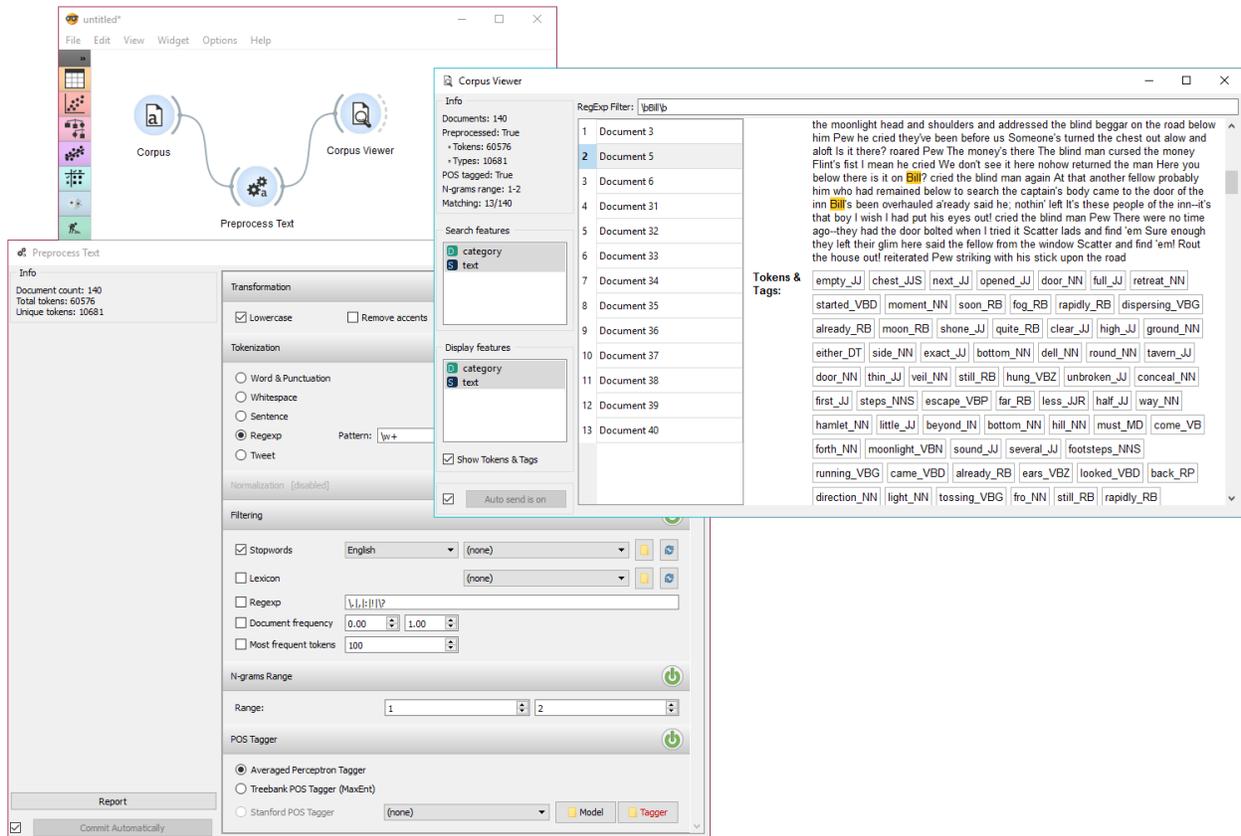


- *POS tagged*: if POS tags are on the input, the result is True, else False.
 - *N-grams range*: if N-grams are set in *Preprocess Text*, results are reported, default is 1-1 (one-grams).
 - *Matching*: number of documents matching the *RegExp Filter*. All documents are output by default.
2. *RegExp Filter*: Python regular expression for filtering documents. By default no documents are filtered (entire corpus is on the output).
 3. *Search Features*: features by which the *RegExp Filter* is filtering. Use Ctrl (Cmd) to select multiple features.
 4. *Display Features*: features that are displayed in the viewer. Use Ctrl (Cmd) to select multiple features.
 5. *Show Tokens & Tags*: if tokens and POS tag are present on the input, you can check this box to display them.
 6. If *Auto commit is on*, changes are communicated automatically. Alternatively press *Commit*.

1.14.1 Example

Corpus Viewer can be used for displaying all or some documents in corpus. In this example, we will first load *book-excerpts.tab*, that already comes with the add-on, into *Corpus* widget. Then we will preprocess the text into words, filter out the stopwords, create bi-grams and add POS tags (more on preprocessing in *Preprocess Text*). Now we want to see the results of preprocessing. In *Corpus Viewer* we can see, how many unique tokens we got and what they are (tick *Show Tokens & Tags*). Since we used also POS tagger to show part-of-speech labels, they will be displayed alongside tokens underneath the text.

Now we will filter out just the documents talking about a character Bill. We use regular expression `\bBill\b` to find the documents containing only the word Bill. You can output matching or non-matching documents, view them in another *Corpus Viewer* or further analyse them.



1.15 Word Cloud

Generates a word cloud from corpus.

Inputs

- Topic** Selected topic.
- Corpus** A collection of documents.

Outputs

- Corpus** Documents that match the selection.
- Word** Selected word that can be used as query in *Concordance*.

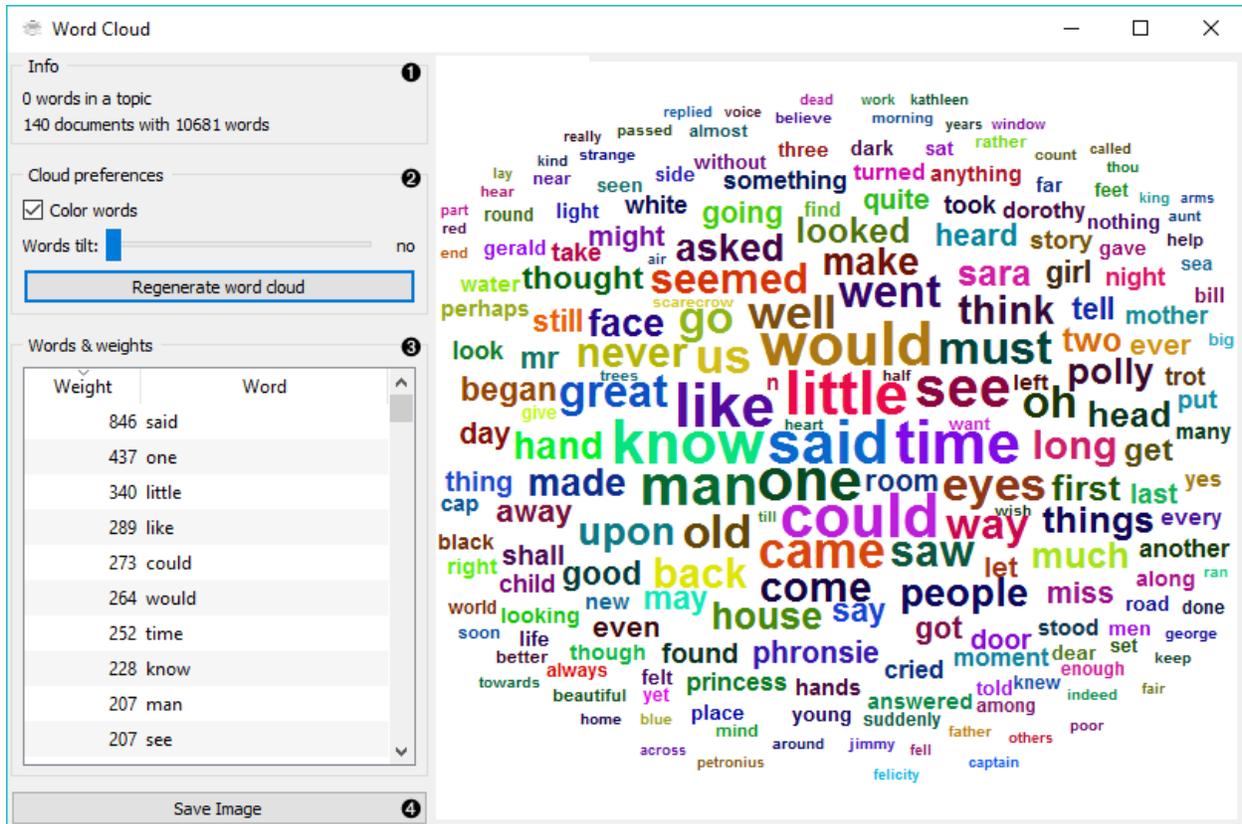
Word Cloud displays tokens in the corpus, their size denoting the frequency of the word in corpus. Words are listed by their frequency (weight) in the widget. The widget outputs documents, containing selected tokens from the word cloud.

1. Information on the input.

- number of words (tokens) in a topic
- number of documents and tokens in the corpus

2. Adjust the plot.

- If *Color words* is ticked, words will be assigned a random color. If unchecked, the words will be black.



- *Word tilt* adjust the tilt of words. The current state of tilt is displayed next to the slider ('no' is the default).
 - *Regenerate word cloud* plot the cloud anew.
3. *Words & weights* displays a sorted list of words (tokens) by their frequency in the corpus or topic. Clicking on a word will select that same word in the cloud and output matching documents. Use *Ctrl* to select more than one word. Documents matching ANY of the selected words will be on the output (logical OR).
 4. *Save Image* saves the image to your computer in a .svg or .png format.

1.15.1 Example

Word Cloud is an excellent widget for displaying the current state of the corpus and for monitoring the effects of preprocessing.

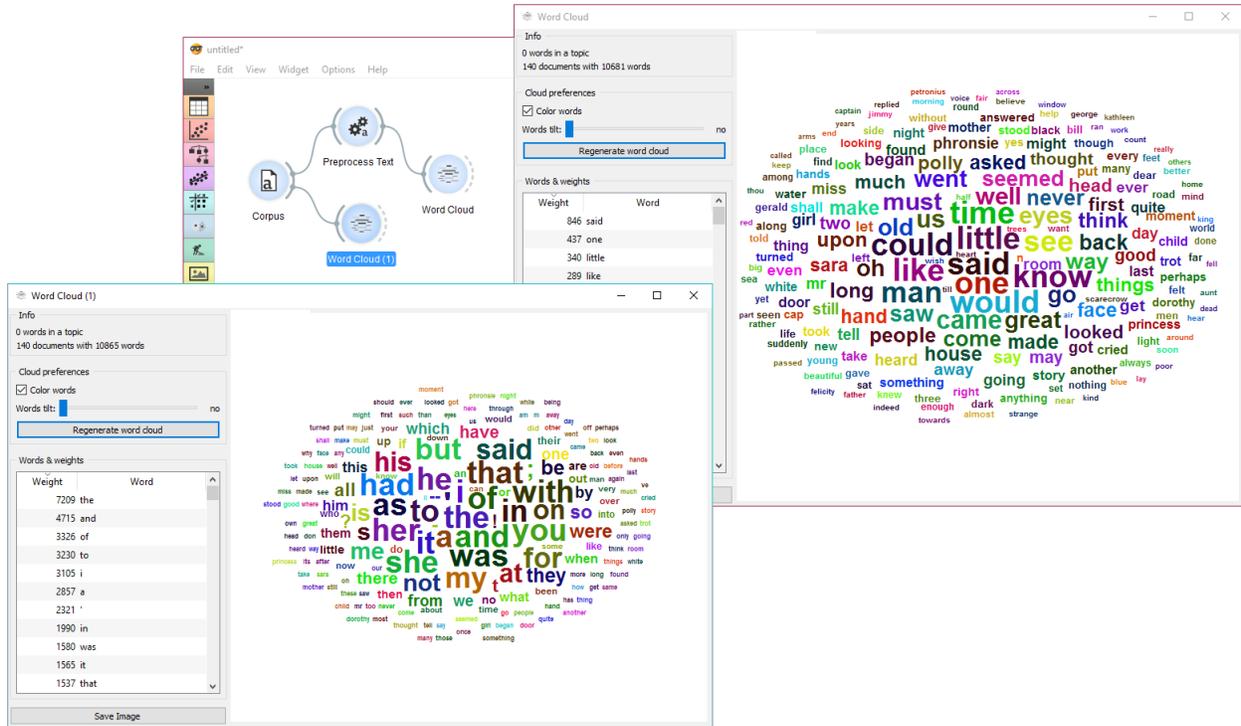
Use *Corpus* to load the data. Connect *Preprocess Text* to it and set your parameters. We've used defaults here, just to see the difference between the default preprocessing in the **Word Cloud** widget and the **Preprocess Text** widget.

We can see from the two widgets, that **Preprocess Text** displays only words, while default preprocessing in the **Word Cloud** tokenizes by word and punctuation.

1.16 Concordance

Display the context of the word.

Inputs



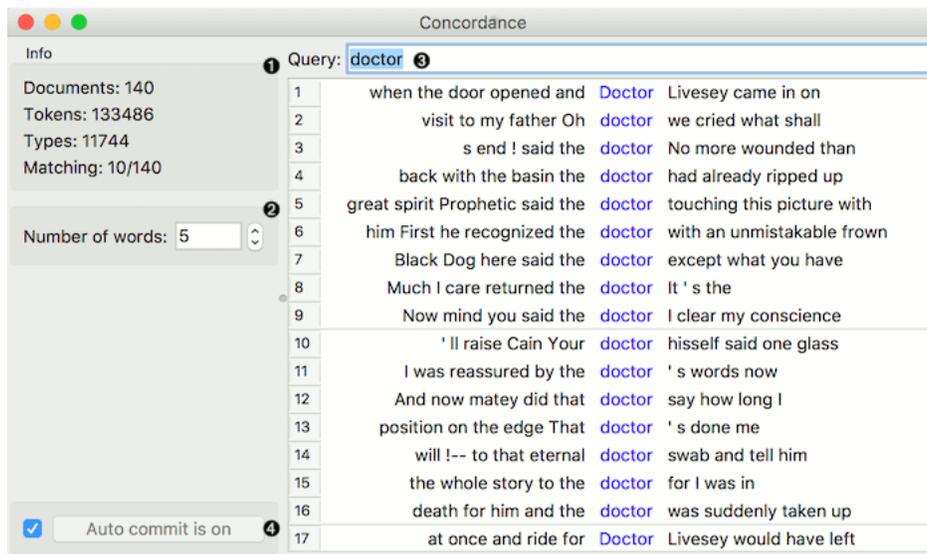
Corpus A collection of documents.

Outputs

Selected Documents Documents containing the queried word.

Concordances A table of concordances.

Concordance finds the queried word in a text and displays the context in which this word is used. Results in a single color come from the same document. The widget can output selected documents for further analysis or a table of concordances for the queried word. Note that the widget finds only exact matches of a word, which means that if you query the word 'do', the word 'doctor' won't appear in the results.



1. Information:

- *Documents*: number of documents on the input.
- *Tokens*: number of tokens on the input.
- *Types*: number of unique tokens on the input.
- *Matching*: number of documents containing the queried word.

2. *Number of words*: the number of words displayed on each side of the queried word.

3. Queried word.

4. If *Auto commit is on*, selected documents are communicated automatically. Alternatively press *Commit*.

1.16.1 Examples

Concordance can be used for displaying word contexts in a corpus. First, we load *book-excerpts.tab* in *Corpus*. Then we connect **Corpus** to **Concordance** and search for concordances of a word “doctor”. The widget displays all documents containing the word “doctor” together with their surrounding (contextual) words.

Now we can select those documents that contain interesting contexts and output them to *Corpus Viewer* to inspect them further.

The screenshot shows the Orange3 interface with the following widgets and data:

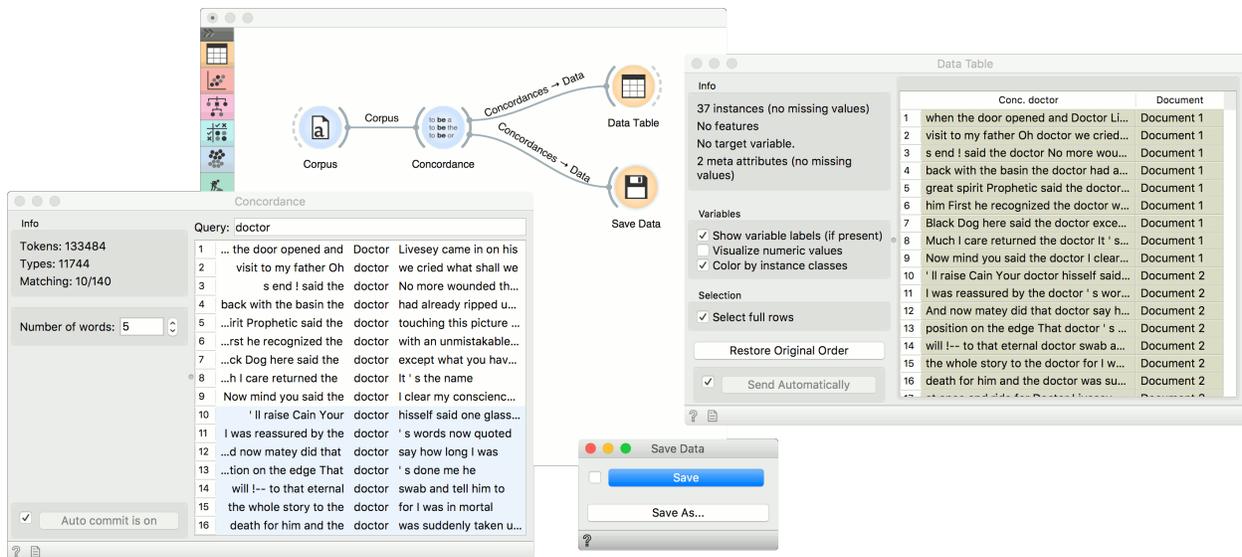
- Corpus**: Info panel shows Documents: 140, Tokens: 133486, Types: 11744, Matching: 112/140. Number of words: 10.
- Concordance**: Query: do. Shows a list of documents with the queried word 'do' highlighted in blue. The list includes:
 - 1 In the meantime we had no idea what to do to help the captain nor any other thought b.
 - 2 ...y father Oh doctor we cried what shall we do ? Where is he wounded ? Wounded ? A
 - 3 ...ossible nothing about it For my part I must do my best to save this fellow ' s trebly
 - 4 ' t break off short you ' ll die -- do you understand that ?-- die and go to your
 - 5 ! he cried A week ! I can ' t do that ; they ' d have the black spot
 - 6 ... ears is singing Lay me back Before I could do much to help him he had fallen back again
 - 7 What I should have done had all gone well I do not know Probably I should have told the w...
 - 8 Bill ' If you don ' t ! ' ll do this and with that he gave me a twitch
 - 9 ...sickness He made a movement to rise but I do not believe he had enough force left in his
 - 10 ' clock I he cried Six hours We ' ll do them yet and he sprang to his feet Even
 - 11 a man would go along with us All they would do was to give me a loaded pistol lest we
 - 12 she gave a sigh and fell on my shoulder I do not know how I found the strength to do
 - 13 I do not know how I found the strength to do it at all and I am afraid it was
 - 14 her for the bridge was too low to let me do more than crawl below it So there we had
 - 15 like a fish out of water and all he could do was to dispatch a man to B ---- to
 - 16 ... were these villains after but money ? What do they care for but money ? For what would
 - 17 Bristol March 1 17 -- Dear Livesey -- As I do not know whether you are at the hall or
- Corpus Viewer**: Shows the full text of the selected document, with the word 'do' highlighted in blue. The text includes:
 - category: children
 - text: the house Jim says he rum ; and as he spoke he reeled a little and caught himself with one hand against the wall Are you hurt? cried I Rum he repeated I must get away from here Rum! Rum! I ran to fetch it but I was quite unsteadied by all that had fallen out and I broke one glass and fouled the tap and while I was still getting in my own way I heard a loud fall in the parlour and running in beheld the captain lying full length upon the floor At the same instant my mother alarmed by the cries and fighting came running downstairs to help me Between us we raised his head He was breathing very loud and hard but his eyes were closed and his face a horrible colour Dear deary me cried my mother what a disgrace upon the house! And your poor father sick! In the meantime we had no idea what to do to help the captain nor any other thought but that he had got his death-wurt in the scuffle with the stranger I got the rum to be sure and tried to put it down his throat but his teeth were tightly shut and his jaws as strong as iron It was a happy relief for us when the door opened and Doctor Livesey came in on his visit to my father On doctor we cried what shall we do? Where is he wounded? Wounded? A fiddle-stick's end! said the doctor No more wounded than you or I The man has had a stroke as I warned him Now Mrs Hawkins just you run upstairs to your husband and tell him if possible nothing about it For my part I must do my best to save this fellow's trebly worthless life; Jim you get me a basin When I got back with the basin the doctor had already ripped up the captain's sleeve and exposed his great sinewy arm It was tattooed in several places Here's luck A fair wind and Billy Bones his fancy were very neatly and clearly executed on the forearm; and up near the shoulder there was a sketch of a gallows and a man hanging from it--done as I thought with great spirit Prophetic said the

In the second example, we will output concordances instead. We will keep the *book-excerpts.tab* in *Corpus* and the connection to **Concordance**. Our queried word remain “doctor”.

This time, we will connect **Data Table** to **Concordance** and select Concordances output instead. In the **Data Table**, we get a list of concordances for the queried word and the corresponding documents. Now, we will save this table with **Save Data** widget, so we can use it in other projects or for further analysis.

1.17 GeoMap

Displays geographic distribution of data.



Inputs

Data Data set.

Outputs

Corpus Documents containing mentions of selected geographical regions.

GeoMap widget shows geolocations from textual (string) data. It finds mentions of geographic names (countries and capitals) and displays distributions (frequency of mentions) of these names on a map. It works with any Orange widget that outputs a data table and that contains at least one string attribute. The widget outputs selected data instances, that is all documents containing mentions of a selected country (or countries).

1. Select the meta attribute you want to search geolocations by. The widget will find all mentions of geolocations in a text and display distributions on a map.
2. Select the type of map you wish to display. The options are *World*, *Europe* and *USA*. You can zoom in and out of the map by pressing + and - buttons on a map or by mouse scroll.
3. The legend for the geographic distribution of data. Countries with the boldest color are most often mentioned in the selected region attribute (highest frequency).

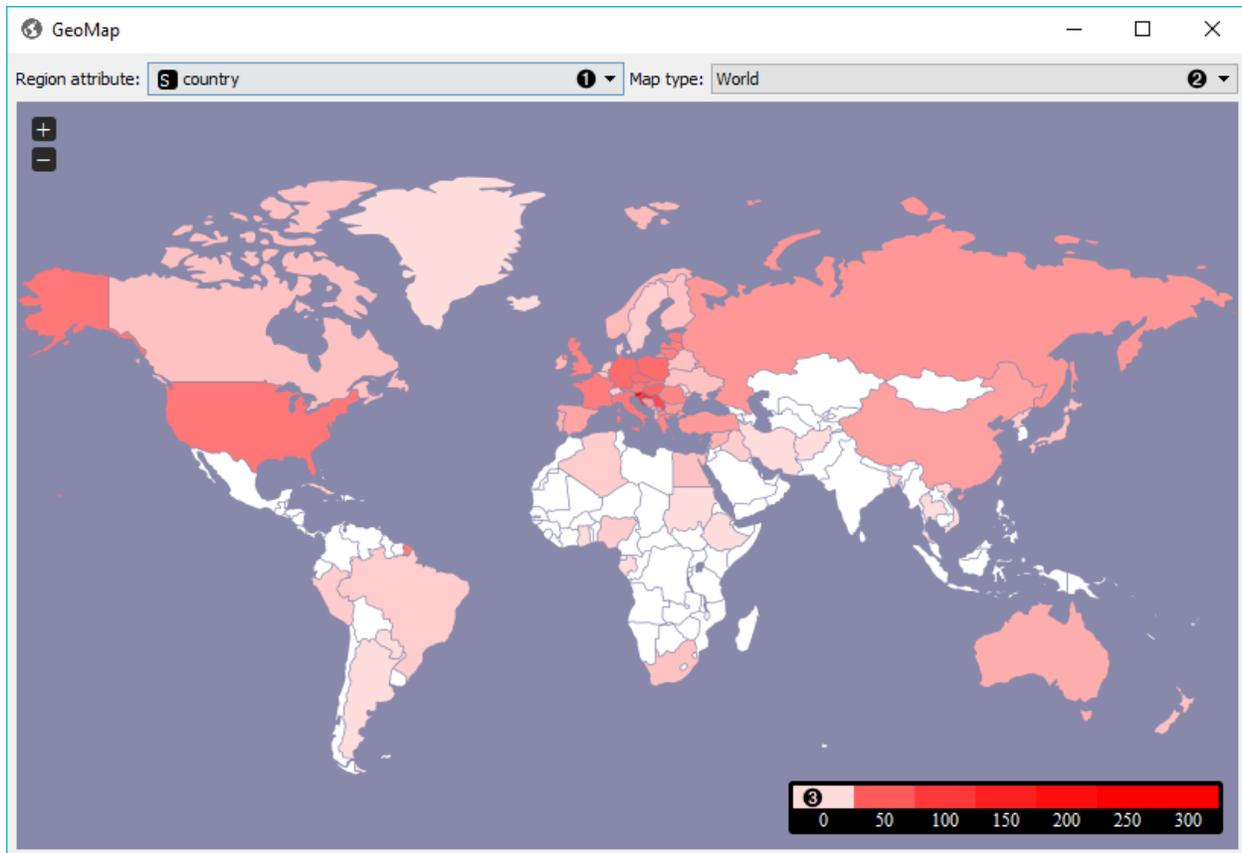
To select documents mentioning a specific country, click on a country and the widget will output matching documents. To select more than one country hold Ctrl/Cmd upon selection.

1.17.1 Example

GeoMap widget can be used for simply visualizing distributions of geolocations or for a more complex interactive data analysis. Here, we've queried *NY Times* for articles on Slovenia for the time period of the last year (2015-2016). First we checked the results with *Corpus Viewer*.

Then we sent the data to **GeoMap** to see distributions of geolocations by *country* attribute. The attribute already contains country tags for each article, which is why *NY Times* is great in combinations with **GeoMap**. We selected Germany, which sends all the documents tagged with Germany to the output. Remember, we queried *NY Times* for articles on Slovenia.

We can again inspect the output with **Corpus Viewer**. But there's a more interesting way of visualizing the data. We've sent selected documents to *Preprocess Text*, where we've tokenized text to words and removed stopwords.



Finally, we can inspect the top words appearing in last year's documents on Slovenia and mentioning also Germany with *Word Cloud*.

1.18 Word Enrichment

Word enrichment analysis for selected documents.

Inputs

Corpus A collection of documents.

Selected Data Selected instances from corpus.

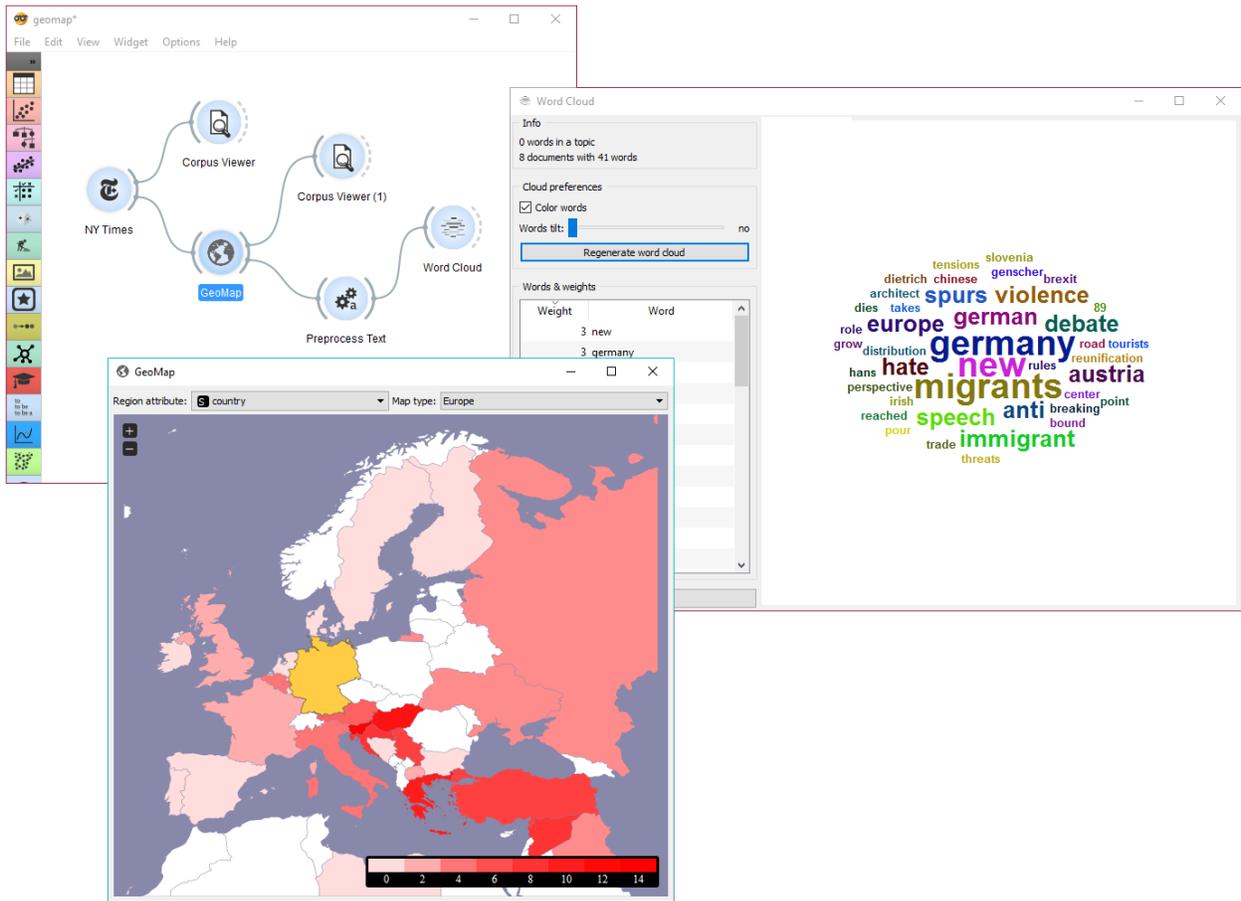
Outputs None

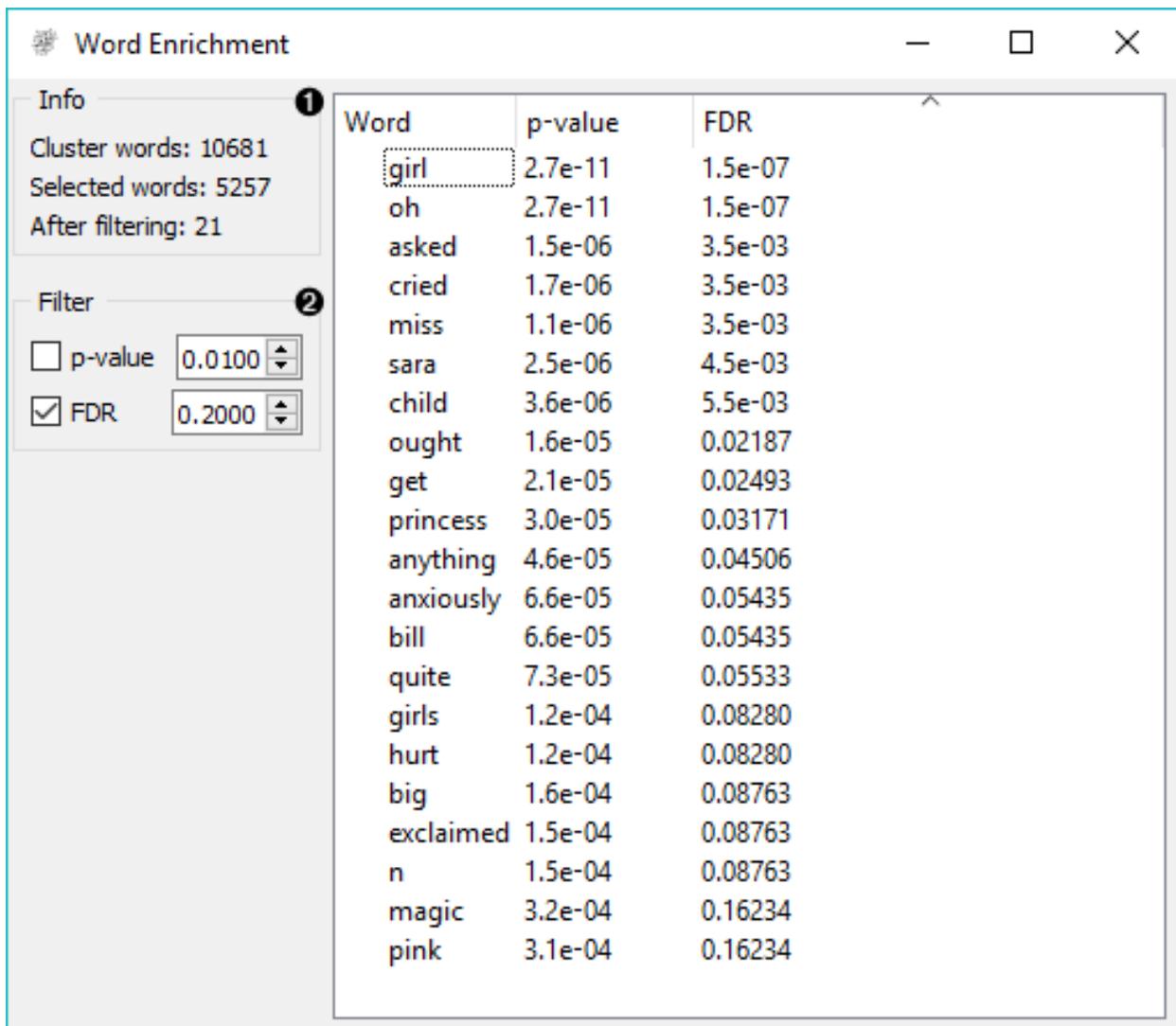
Word Enrichment displays a list of words with lower p-values (higher significance) for a selected subset compared to the entire corpus. Lower p-value indicates a higher likelihood that the word is significant for the selected subset (not randomly occurring in a text). FDR (False Discovery Rate) is linked to p-value and reports on the expected percent of false predictions in the set of predictions, meaning it account for false positives in list of low p-values.

1. Information on the input.

- Cluster words are all the tokens from the corpus.
- Selected words are all the tokens from the selected subset.
- After filtering reports on the enriched words found in the subset.

2. Filter enables you to filter by:





Word Enrichment

Info ①

Cluster words: 10681
Selected words: 5257
After filtering: 21

Filter ②

p-value 0.0100

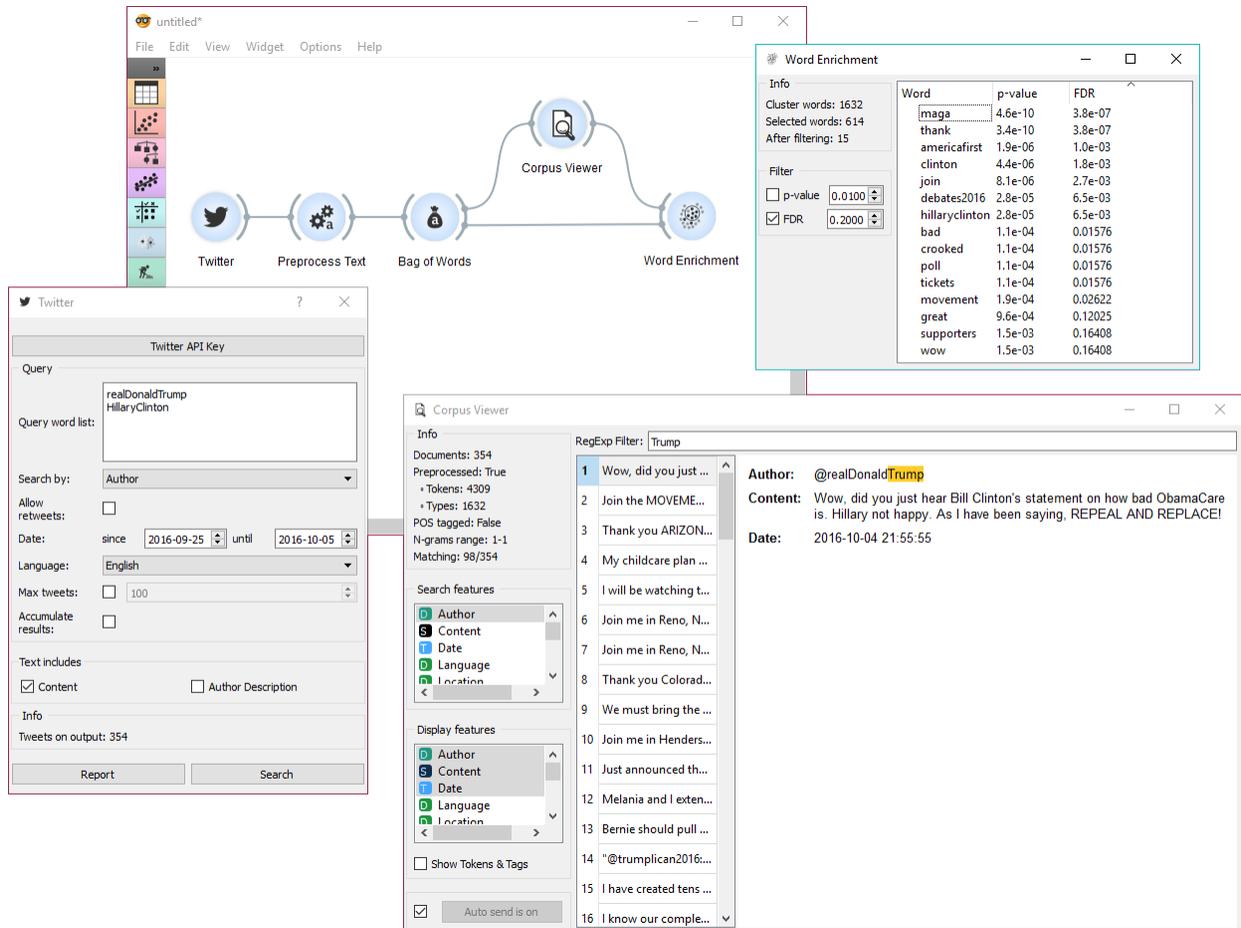
FDR 0.2000

Word	p-value	FDR
girl	2.7e-11	1.5e-07
oh	2.7e-11	1.5e-07
asked	1.5e-06	3.5e-03
cried	1.7e-06	3.5e-03
miss	1.1e-06	3.5e-03
sara	2.5e-06	4.5e-03
child	3.6e-06	5.5e-03
ought	1.6e-05	0.02187
get	2.1e-05	0.02493
princess	3.0e-05	0.03171
anything	4.6e-05	0.04506
anxiously	6.6e-05	0.05435
bill	6.6e-05	0.05435
quite	7.3e-05	0.05533
girls	1.2e-04	0.08280
hurt	1.2e-04	0.08280
big	1.6e-04	0.08763
exclaimed	1.5e-04	0.08763
n	1.5e-04	0.08763
magic	3.2e-04	0.16234
pink	3.1e-04	0.16234

- p-value
- false discovery rate (FDR)

1.18.1 Example

In the example below, we're retrieved recent tweets from the 2016 presidential candidates, Donald Trump and Hillary Clinton. Then we've preprocessed the tweets to get only words as tokens and to remove the stopwords. We've connected the preprocessed corpus to *Bag of Words* to get a table with word counts for our corpus.



Then we've connected *Corpus Viewer* to **Bag of Words** and selected only those tweets that were published by Donald Trump. See how we marked only the *Author* as our *Search feature* to retrieve those tweets.

Word Enrichment accepts two inputs - the entire corpus to serve as a reference and a selected subset from the corpus to do the enrichment on. First connect **Corpus Viewer** to **Word Enrichment** (input Matching Docs → Selected Data) and then connect **Bag of Words** to it (input Corpus → Data). In the **Word Enrichment** widget we can see the list of words that are more significant for Donald Trump than they are for Hillary Clinton.

1.19 Duplicate Detection

Detect & remove duplicates from a corpus.

Inputs

Distances A distance matrix.

Outputs

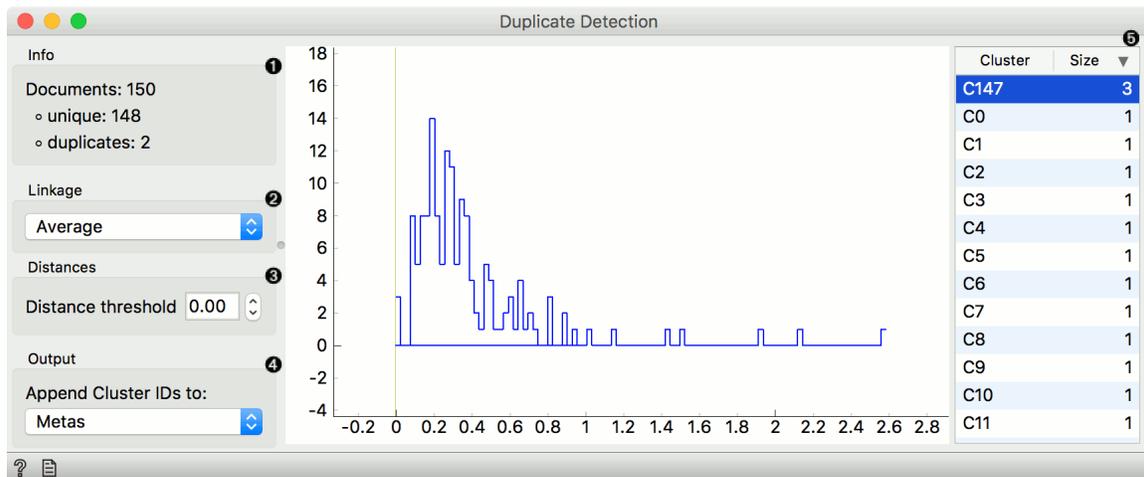
Corpus Without Duplicated Corpus with duplicates removed.

Duplicates Cluster Documents belonging to selected cluster.

Corpus Corpus with appended cluster labels.

Duplicate Detection uses clustering to find duplicates in the corpus. It is great with the *Twitter* widget for removing retweets and other similar documents.

To set the level of similarity, drag the line vertical line left or right in the visualization. The further left the line, the more similar the documents have to be in order to be considered duplicates. You can also set the threshold manually in the control area.



1. Information on unique and duplicate documents.
2. Linkage used for clustering (Single, Average, Complete, Weighted and Ward).
3. Distance threshold sets the similarity cutoff. The lower the value, the more similar the data instances have to be to belong to the same cluster. You can also set the cutoff by dragging the vertical line in the plot.
4. Cluster labels can be appended as attributes, class or metas.
5. List of clusters at the selected threshold. They are sorted by size by default. Click on the cluster to observe its content on the output.

1.19.1 Example

This simple example uses *iris* data to find identical data instances. Load *iris* with the **File** widget and pass it to **Distances**. In **Distances**, use Euclidean distance for computing the distance matrix. Pass distances to **Duplicate Detection**.

It looks like cluster C147 contain three duplicate entries. Let us select it in the widget and observe it in a **Data Table**. Remember to set the output to *Duplicates Cluster*. The three data instances are identical. To use the data set without duplicates, use the first output, *Corpus Without Duplicates*.

The same procedure can be used also for corpora. Remember to use the *Bag of Words* between *Corpus* and **Distances**.

The screenshot displays an Orange3 workflow and several widget windows. The workflow consists of the following widgets: File, Distances, Duplicate Detection, and Data Table. The Duplicate Detection widget is currently open, showing the following settings and results:

Distances widget settings:

- Distances between: Rows
- Distance Metric: Euclidean
- Normalized:
- Apply Automatically:

Duplicate Detection widget info:

- Documents: 150
 - unique: 148
 - duplicates: 2
- Linkage: Average
- Distances: Distance threshold 0.00
- Output: Append Cluster IDs to: Metas

The Duplicate Detection widget also displays a dendrogram and a cluster list:

Cluster	Size
C147	3
C0	1
C1	1
C2	1
C3	1
C4	1
C5	1
C6	1
C7	1
C8	1
C9	1
C10	1
C11	1

The Data Table widget shows the following data:

	iris	sepal length	sepal width	petal length	petal width
1	Iris-setosa	4.9	3.1	1.5	0.1
2	Iris-setosa	4.9	3.1	1.5	0.1
3	Iris-setosa	4.9	3.1	1.5	0.1

2.1 Corpus

2.2 Preprocessor

2.3 Twitter

2.4 New York Times

2.5 The Guardian

2.6 Wikipedia

2.7 Bag of Words

2.8 Topic Modeling

2.9 Tag

CHAPTER 3

Indices and tables

- `genindex`
- `modindex`
- `search`