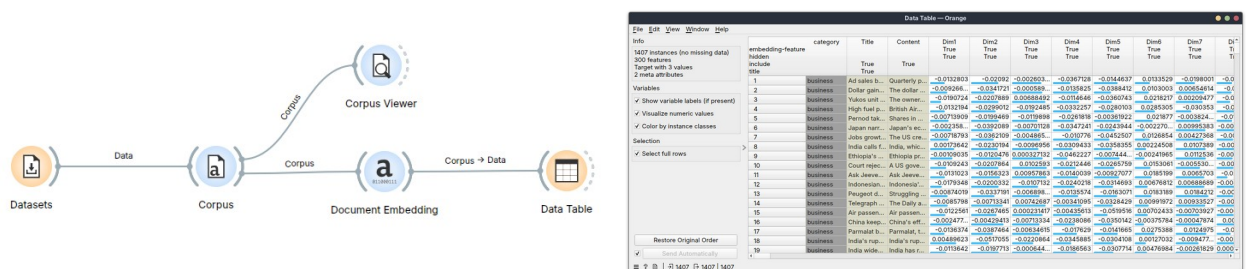# No code lab

## Text processing

*Hint: before doing this exercise you have to install Orange-Text extensions. Open the [Options] menu in the main window, then choose [Add-ons], search for Text and click the checkbox to mark the add-on for installation. Click OK. You will most probably have to restart the tool after add-on installation.*

1. Use **Datasets** operator to load the *BBC3* dataset. This dataset contains over 1400 news articles (titles and contents) classified into several classes based on article topics (sports, politics, entertainment, etc.)
2. Send the data to the **Corpus** operator from the [Text Mining] section. Double-click the operator and make sure that *Used text features* is Content, and *Ignored text features* is Title. Send the output of the **Corpus** operator to the **Corpus Viewer** operator from the [Text Mining] section. Double-click the **Corpus Viewer** to see the articles.
3. First, we will compute the text embedding. Send the output of the **Corpus** operator to the **Document Embedding** operator from the [Text Mining] section. Double-click the **Document Embedding** operator to see the options. Orange Data Mining allows you to use two embeddings, a Sentence BERT (SBERT) and fastText embeddings. Select fastText, make sure that you are using English language embeddings, and use "mean" as the aggregator. FastText computes the embedding for each word, so the embedding of the article will be the mean of all text embeddings. Send the output of the **Document Embedding** operator to the **Data Table** and look at the embeddings. Now you have a very good representation of the text data and you could use this representation to do document clustering and document classification.



4. We will return to a much simpler text representation which is a bag-of-words model. Send the output of the **Corpus** operator to the **Preprocess Text** operator from the [Text Mining] section. Double-click on this operator, as you can see, a lot of operations can be performed here. Make sure that you tranform all words to lowercase. Tokenization of the text can be performed by Word_Punctuation. Use stop-word filtering (the operator will remove very popular and common words such as "the", "in", "for", etc.) To speed up computations check the *Most frequent tokens* checkbox and leave only 1000 most frequent words. Send the output of the **Preprocess Text** operator to the **Word Cloud** operator from the [Text Mining] and see the results.

5. Send the output of the **Corpus** operator to the **Bag of Words** operator from the [Text Mining] section. This operator transforms each article to a vector, where each position represents a word and the number represents the number of occurrences of that word in the article. If you want, you can use the **Data Table** operator to look at the **Bag of Words** transformation.

6. Send the **Bag of Words** output to the **Distances** operator from the [Unsupervised] section. This operator will compute the similarity between articles based on the number of co-occurrences of words.

7. Next, send the computed distances to the **Hierarchical Clustering** operator from the [Unsupervised] section. Finally, send the output of the **Hierarchical Clustering** operator to a **Corpus Viewer** operator. This operator will always display documents that are in the currently selected cluster. Your final workflow should be as follows:



8. Double-click the **Hierarchical Clustering** operator and select one of the clusters. Examine the contents of the cluster in the **Corpus Viewer**. Do all articles from the selected cluster share the same category?

## Assignment

Repeat steps 6-8 for the document representation computed from word embeddings. Compute the distances between embeddings, run hierarchical clustering and examine the clusters of documents.