**Data Science**

# Introduction to Orange3

# Training Series

## Week 3

**2 Februari 2024**

Orange3
Prepared by I Made Murwantara, Universitas Pelita Harapan

# Overview

1. Image Clustering
2. Text Analytics
3. Deployment in Streamlit

Add-on:
- Image Analytics
- Text Mining

# Image Embedding

- a lower-dimensional representation of the image

- a **dense vector representation of the image** which can be used for many tasks such as classification

- Embeddings are different from images in their raw form. An image file contains RGB data that says exactly what colour each pixel is.

- Embeddings encode information that represents the contents of an image.

- These embeddings are unintelligible in their raw form, just as images are when read as a list of numbers.

- It is when you use embeddings that they start to make sense.

https://www.activeloop.ai/resources/generate-image-embeddings-using-a-pre-trained-cnn-and-store-them-in-hub/
https://blog.roboflow.com/what-is-an-image-embedding/

# Image Embedding



- This image contains a bowl of fruit.

- An image embedding will encode this information

- We could then compare the image embedding to a text embedding like "fruit" to see how similar the concept of "fruit" is to the contents of the image.

- We could take two prompts, such as "fruit" and "vegetable", and see how similar each one is.
- The most similar prompt is considered the most representative of the image.

https://blog.roboflow.com/what-is-an-image-embedding/

# Image Embedding

- Deep learning is used to develop models that transform complex objects to vectors of numbers.

- Deep learning requires a lot of data (thousands, possibly millions of data instances) and processing power to prepare the network.
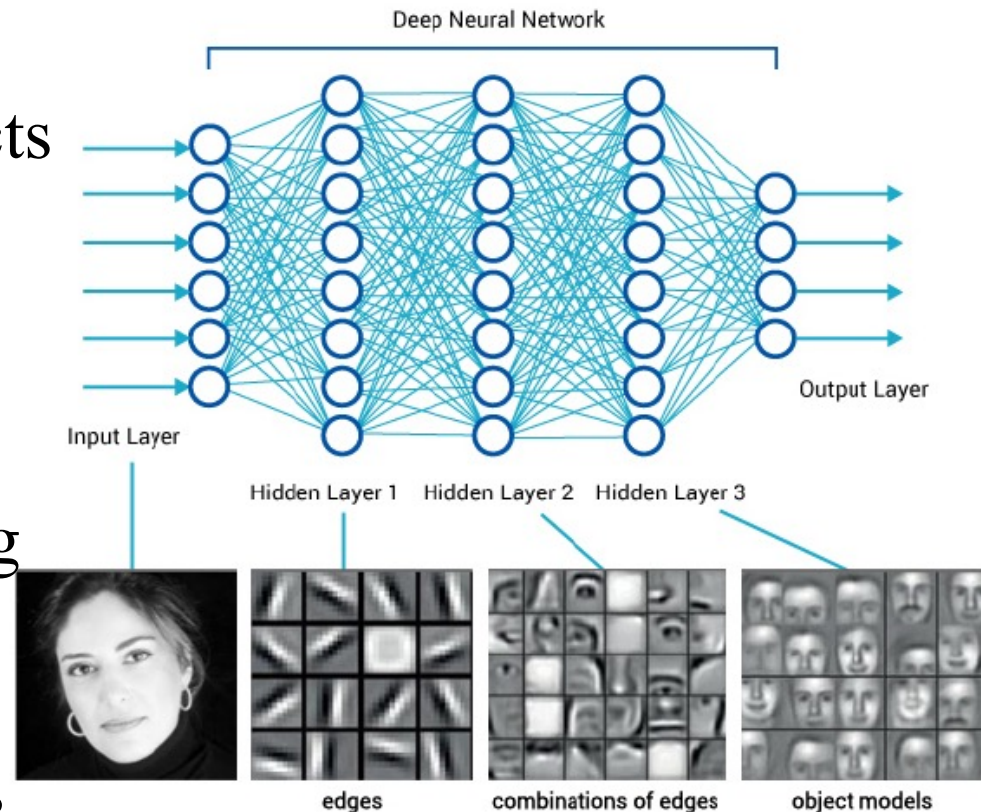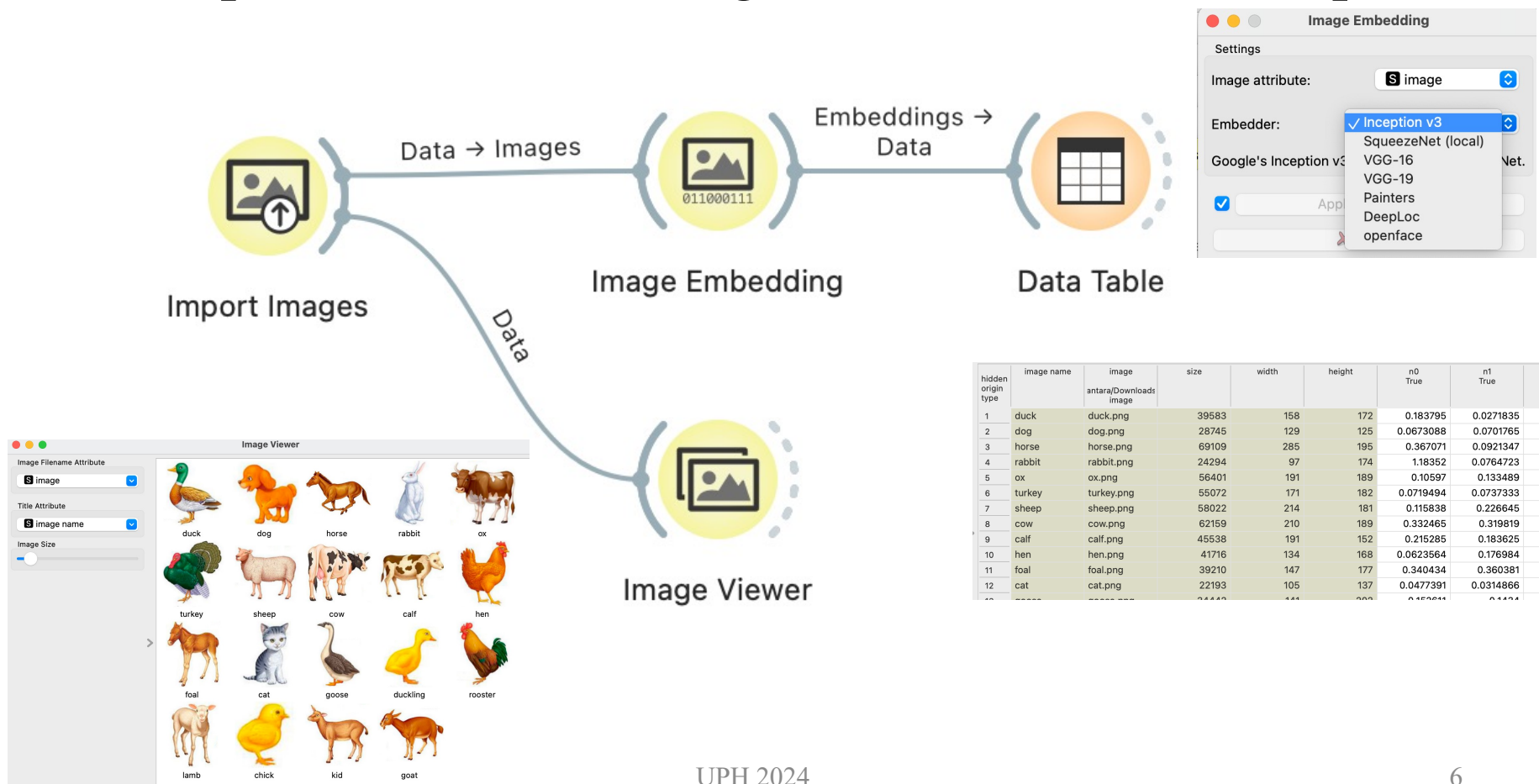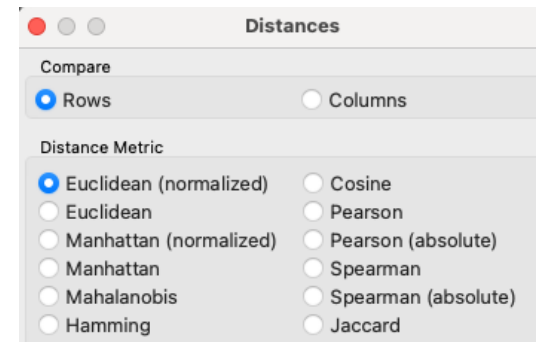
- We will use one which is already prepared.



Deep Neural Network

Input Layer

Hidden Layer 1    Hidden Layer 2    Hidden Layer 3

Output Layer

edges    combinations of edges    object models

# Image Embedding

- Dataset:
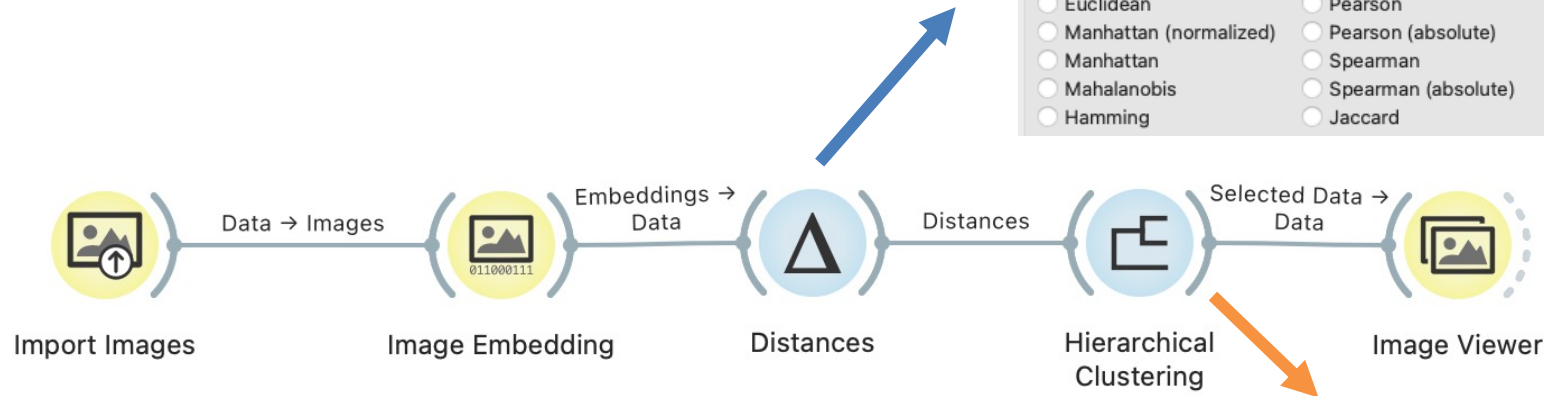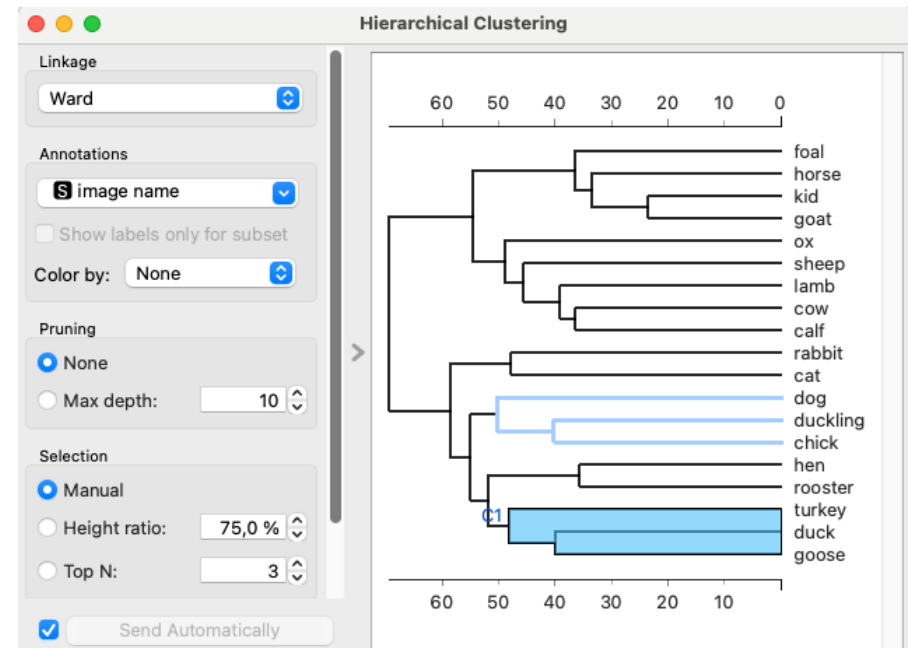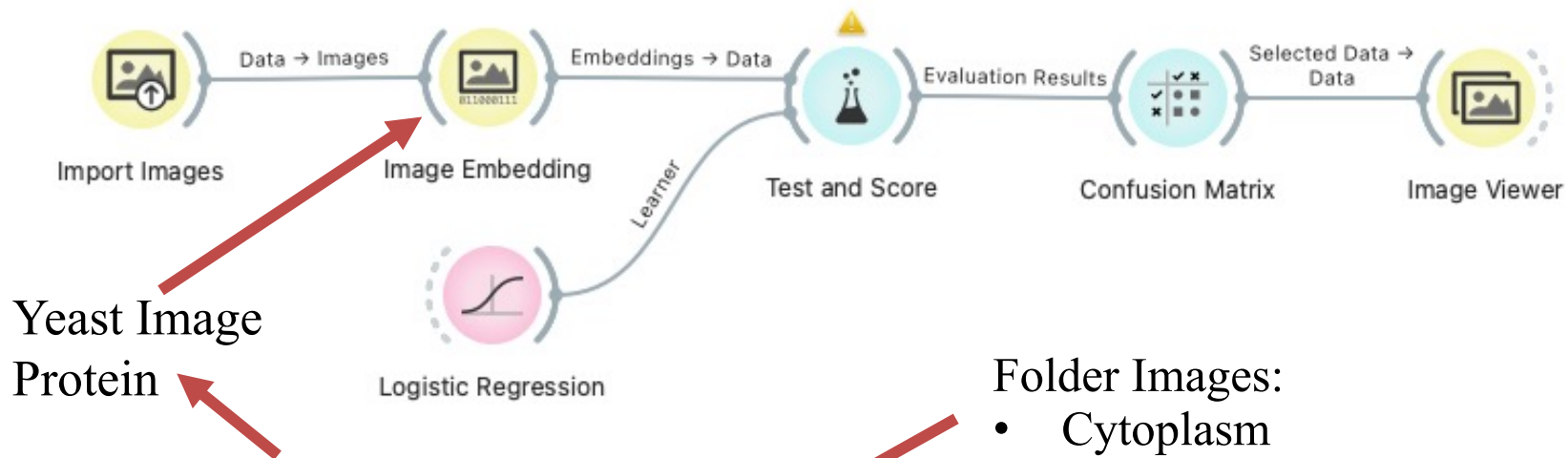  - http://file.biolab.si/images/domestic-animals.zip

# Image Embedding



1. Load images.
2. Turned images into numbers.
3. Distances widget computes distances between rows or columns in a dataset.
4. Group items visualization
5. View group of images

# Images Classification

http://file.biolab.si/files/yeast-localization-small.zip
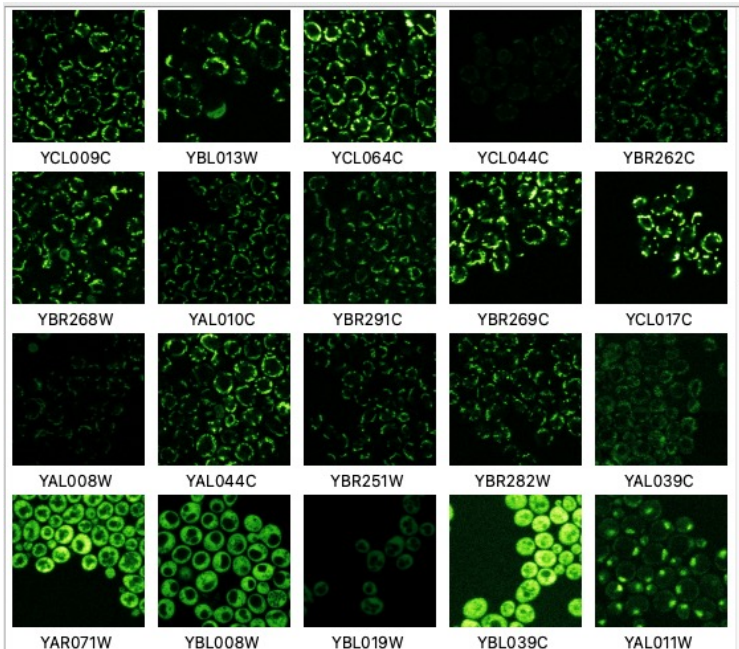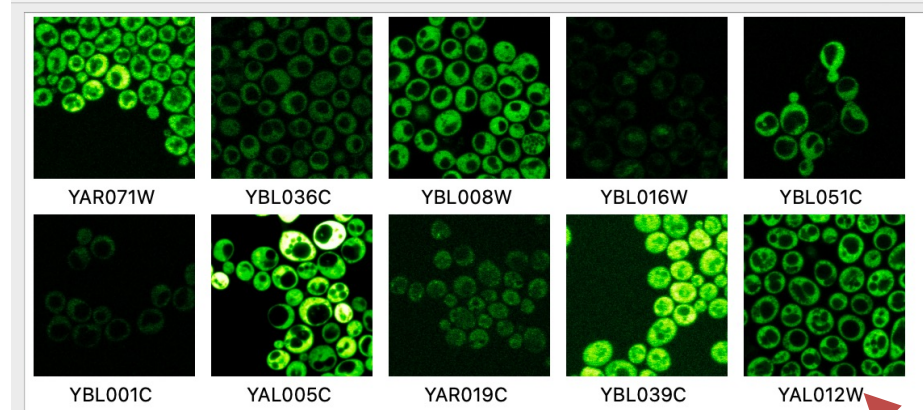


Yeast Image Protein

Folder Images:
- Cytoplasm
- Endosome
- Mitochondria
- Nucleus

Steps
1. Import images
2. Construct matrix table from images
3. Evaluates using Machine Learning
4. Results of Testing Classification
5. View Images

# Images Classification

# Text Mining



Text Mining & Text Analysis – Identifies textual patterns & trends within unstructured data through the use of machine leaning, statistics & linguistics [IBM]
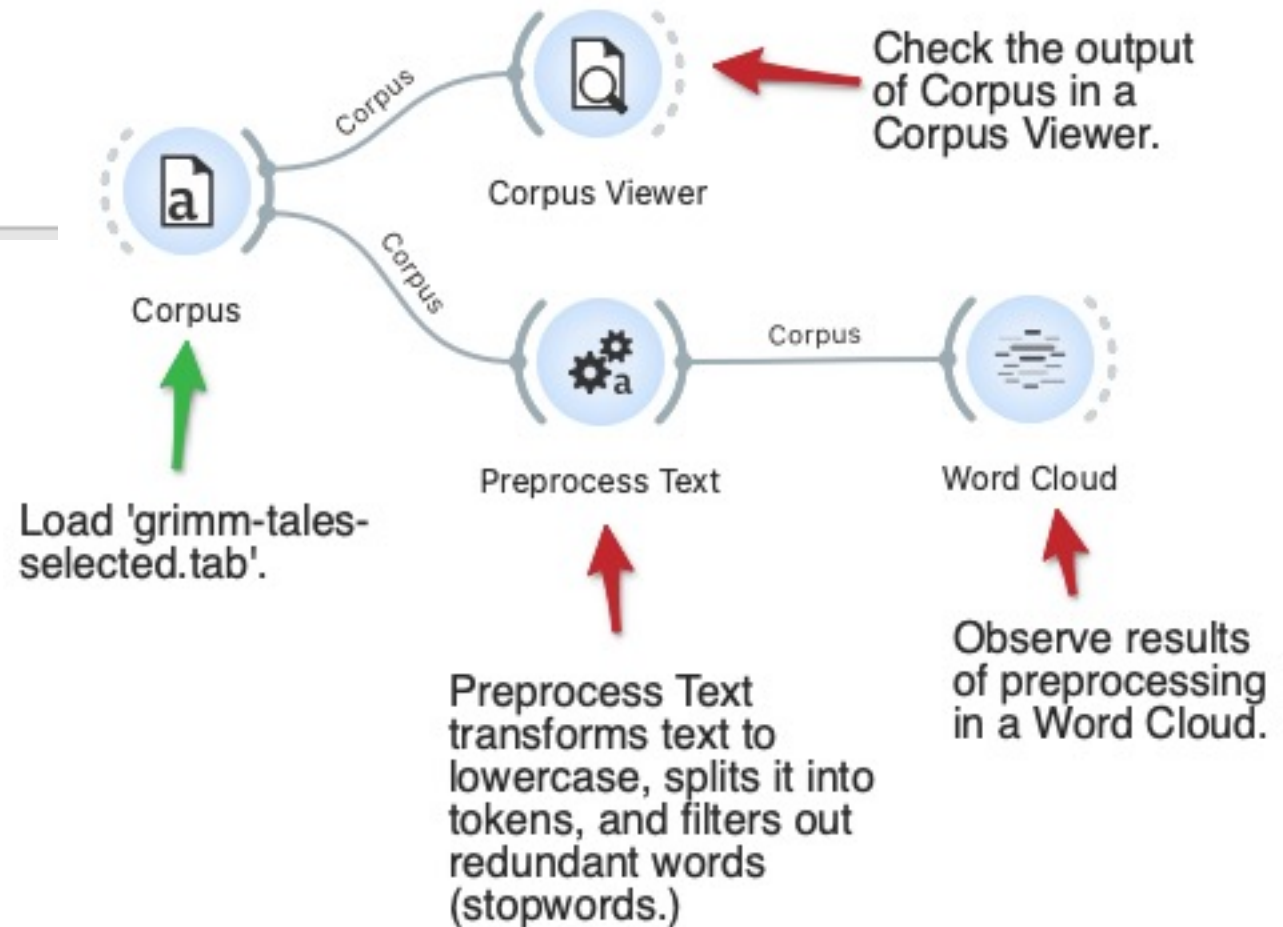
Text Mining is the process of obtaining meaningful information from large collections of unstructured data using Natural Language Processing (NLP)
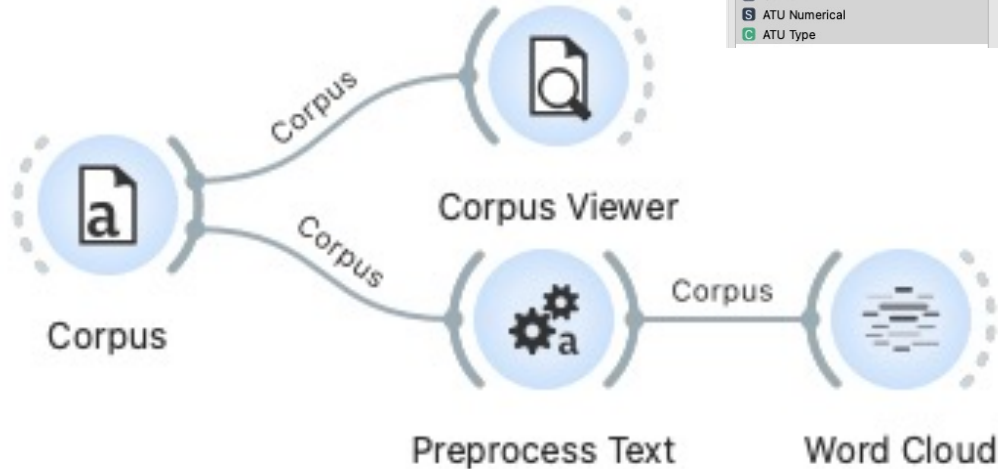
Text mining is the data mining technique or process which discovers earlier unfamiliar and valuable information from a huge quantity of unstructured text data

# Text Preprocessing

**Preprocessors**

- Transformation
- Tokenization
- Normalization
- Filtering
- N-grams Range
- POS Tagger

Check the output of Corpus in a Corpus Viewer.

Corpus Viewer

Corpus

Load 'grimm-tales-selected.tab'.

Preprocess Text

Preprocess Text transforms text to lowercase, splits it into tokens, and filters out redundant words (stopwords.)

Word Cloud

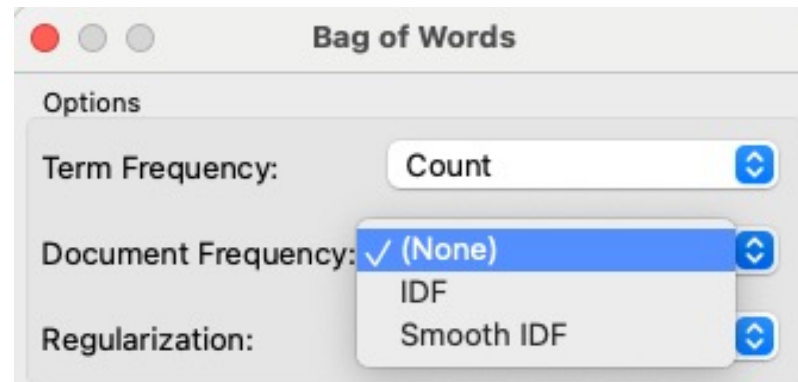Observe results of preprocessing in a Word Cloud.
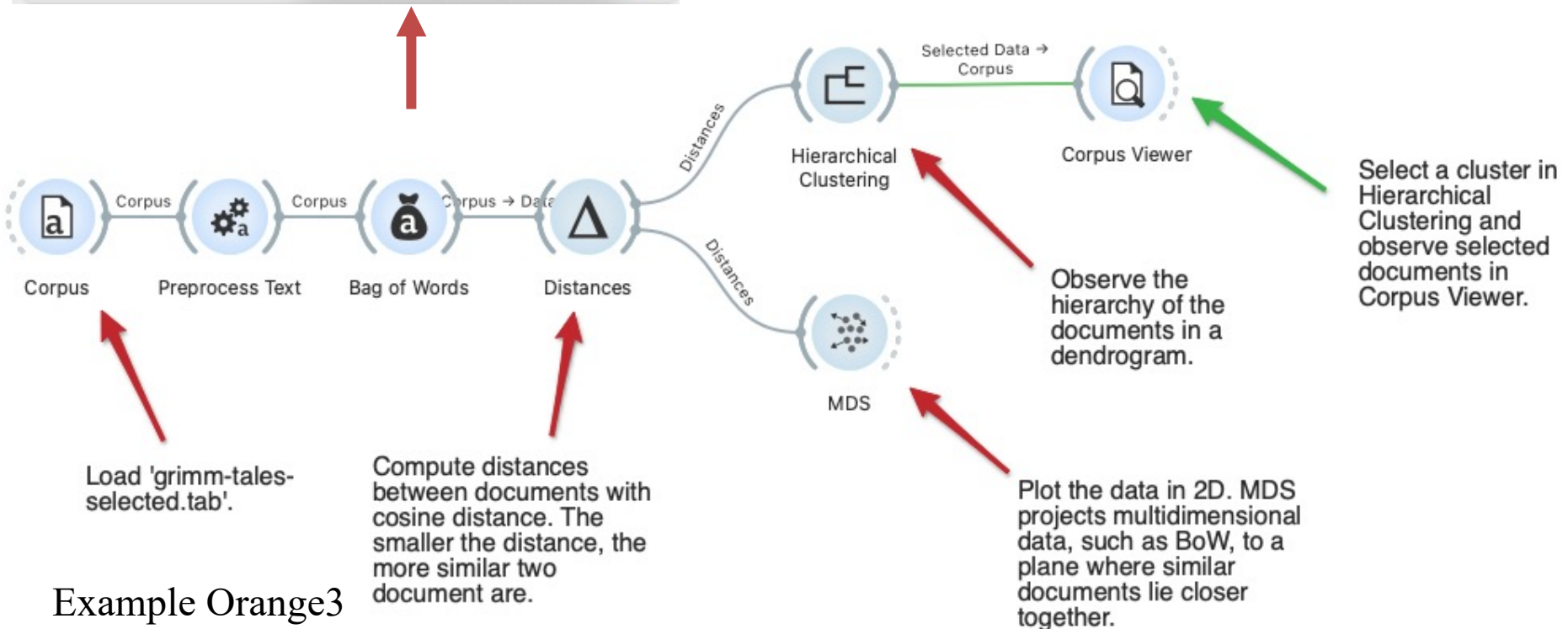
# Text Preprocessing

# Hierarchical Clustering



**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.

Example Orange3

Load 'grimm-tales-selected.tab'.

Compute distances between documents with cosine distance. The smaller the distance, the more similar two document are.

Observe the hierarchy of the documents in a dendrogram.

Select a cluster in Hierarchical Clustering and observe selected documents in Corpus Viewer.

Plot the data in 2D. MDS projects multidimensional data, such as BoW, to a plane where similar documents lie closer together.
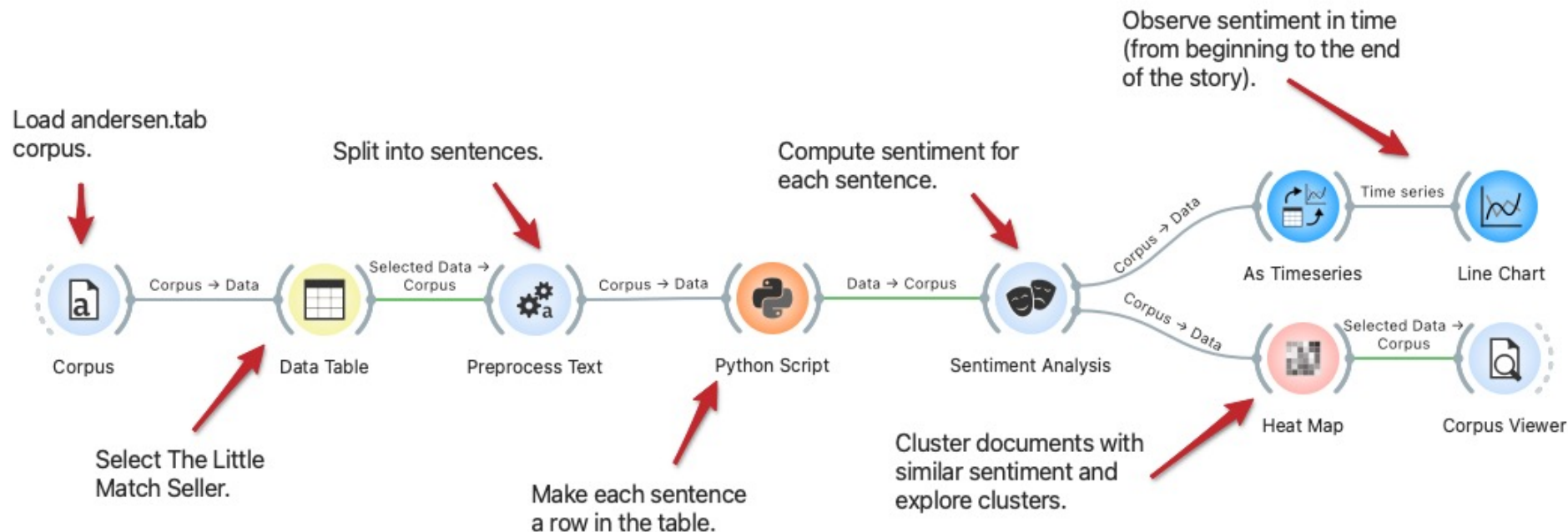
# Text Mining

- Analytics in Text Mining

# Story Arcs

1. Select the story from the corpus of Andersen tales.
2. Create a table where each sentences of the tale into a separate row.
3. Sentiment analysis to compute the sentiment of each sentence, then we observe the emotional arcs through the story
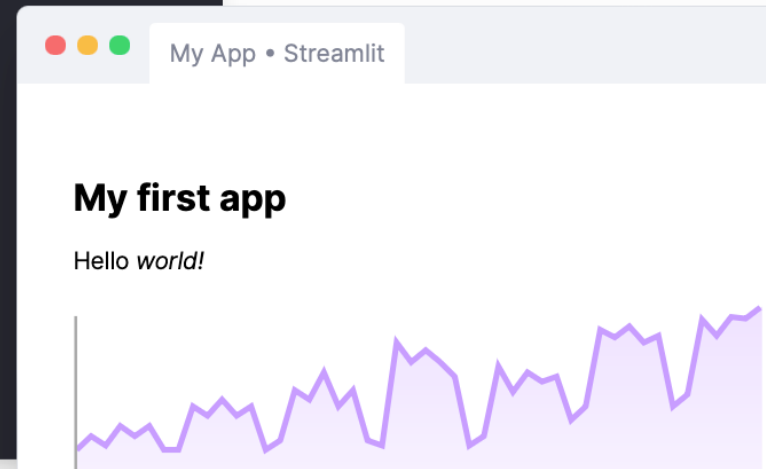4. Observe sentences with similar scores in the Heat Map and Corpus Viewer



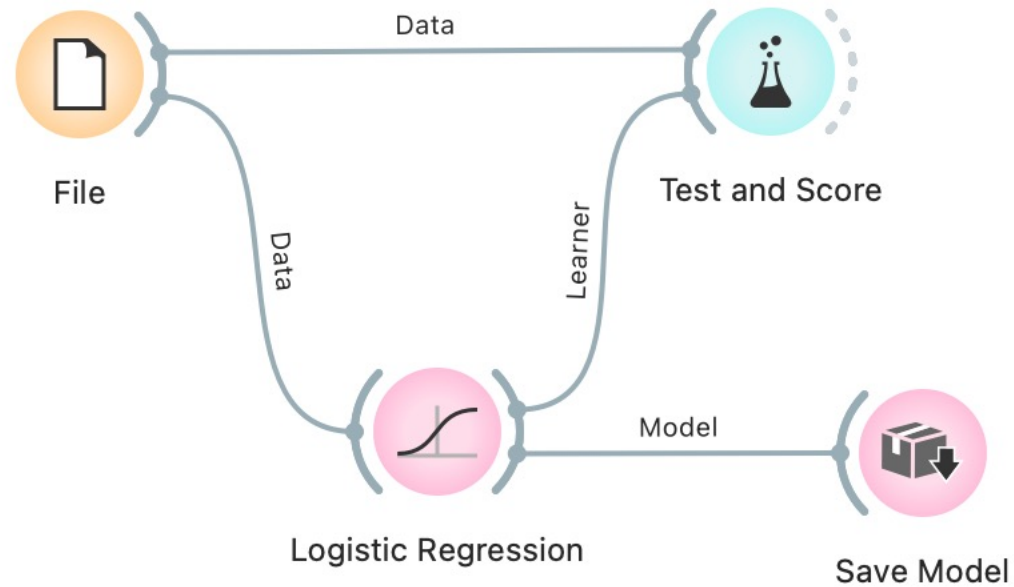Observe sentiment in time (from beginning to the end of the story).

Load andersen.tab corpus.

Split into sentences.

Compute sentiment for each sentence.

Corpus — Corpus → Data — Data Table — Selected Data → Corpus — Preprocess Text — Corpus → Data — Python Script — Data → Corpus — Sentiment Analysis — Corpus → Data — As Timeseries — Time series — Line Chart

Corpus → Data — Heat Map — Selected Data → Corpus — Corpus Viewer

Select The Little Match Seller.

Make each sentence a row in the table.

Cluster documents with similar sentiment and explore clusters.

https://orangedatamining.com/examples/?tag=Text+Mining
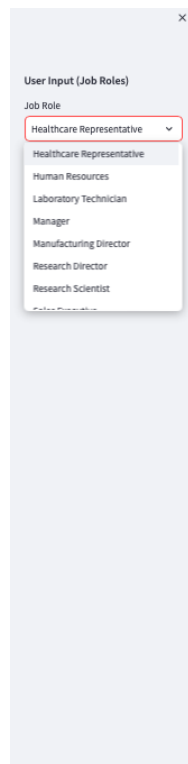
# Deployment

- Web Services
- Python Application
- Open Source

# Deployment

# Deployment

1. Create folder & Copy source files
2. Run "streamlit run files.py"
3. Type ctrl+c to stop

# Any Question ?

- [https://orangedatamining.com/examples/](https://orangedatamining.com/examples/)
- Streamlit.io