# Read map

| 시간 | 내용 |
|---|---|
| 10:00 | Registration |
| 10:30-12:30 | Introduction to Python<br>Name, Namespace, Strings,<br>Functions, List, Tuple, Dictionary |
| 12:30-2:00 | Lunch |
| 2:00-3:30 | File handling/ Class<br>Learning Python Library<br>(matplotlib, pandas) |
| 3:30-5:00 | Simple & Fun<br>Project(Wordcloud and collocation)<br>Online resources |

**01** Introduction to Python and Orange

**02** Machine Learning with Orange

**03** Text Mining with Orange

# Part 1
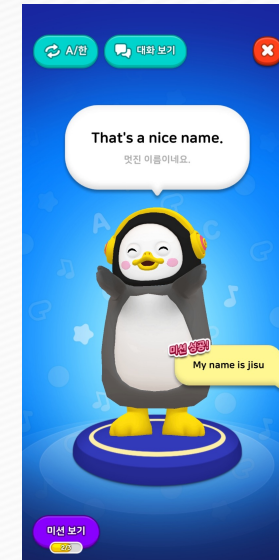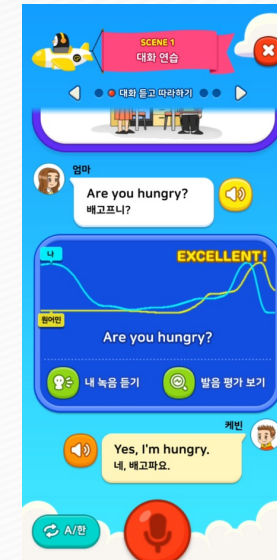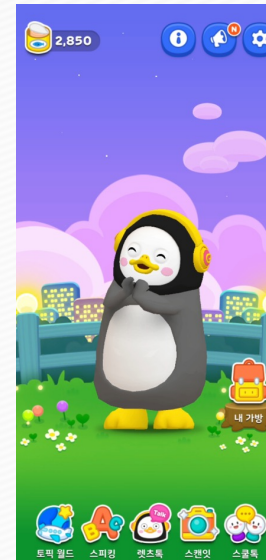# Introduction to Python & Orange

# English Language Teaching & Artificial Intelligence
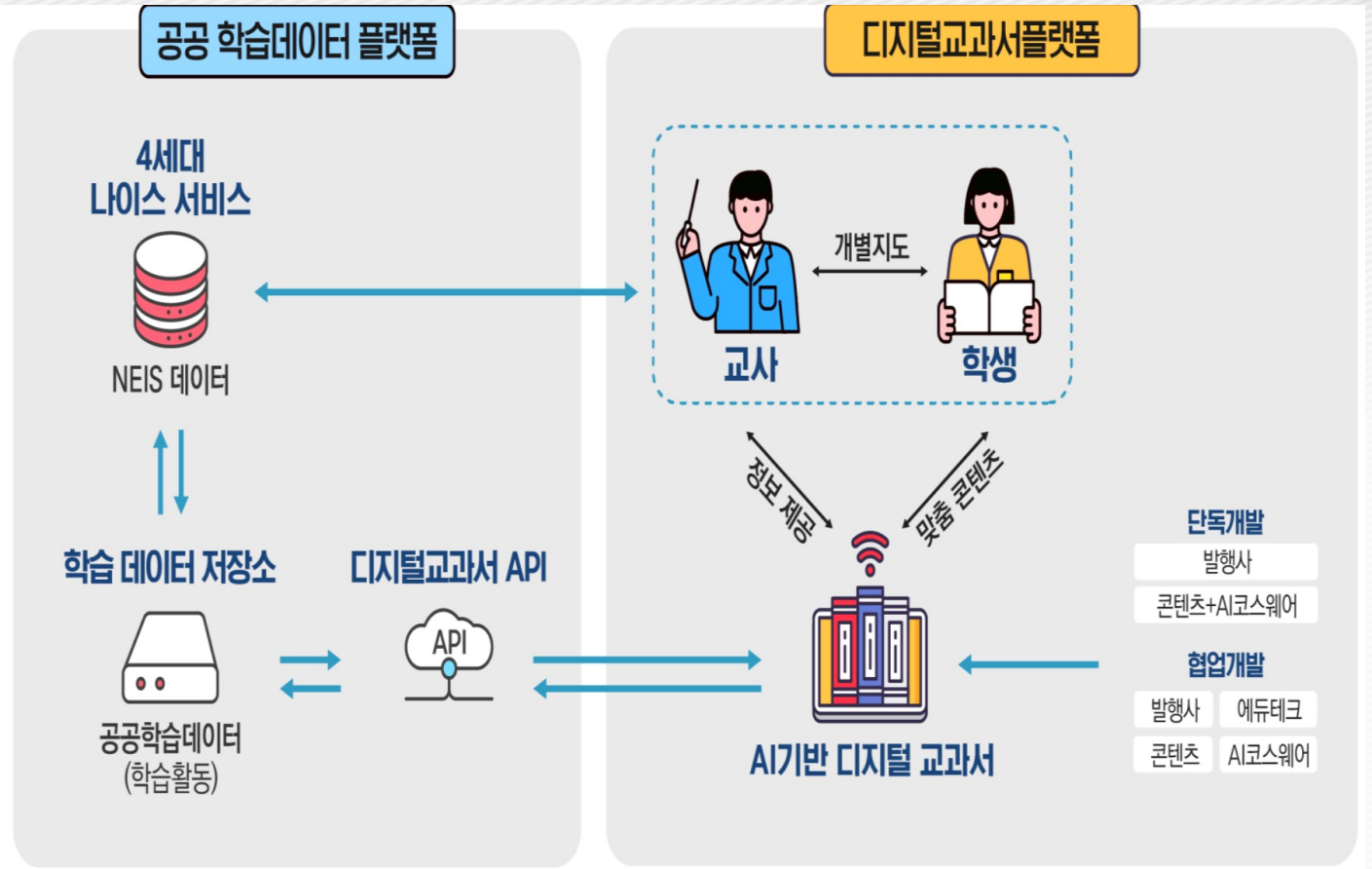
- 2022개정 영어과 교육과정

(바) 다양한 디지털 평가 도구를 적극적으로 활용한다. 디지털 분석·평가 도구를 활용하여 실제적인 평가 맥락을 제공하고 다양한 학습자 데이터를 체계적으로 구축한다. 이를 토대로 다각적이고 신뢰할 만한 평가 결과를 도출할 수 있다.



[그림] 영어과 역량 및 영역 구성

# English Language Teaching & Artificial Intelligence

# GUI vs. CLI



## CLI vs GUI

### Computer Interface definition and explanation

# Guido van Rossum
# 귀도 반 로섬

**89** 크리스마스 연휴

**99** **DARPA**

**Computer Programming for Everybody**

MONTY PYTHON
Film Maker
Lived Here
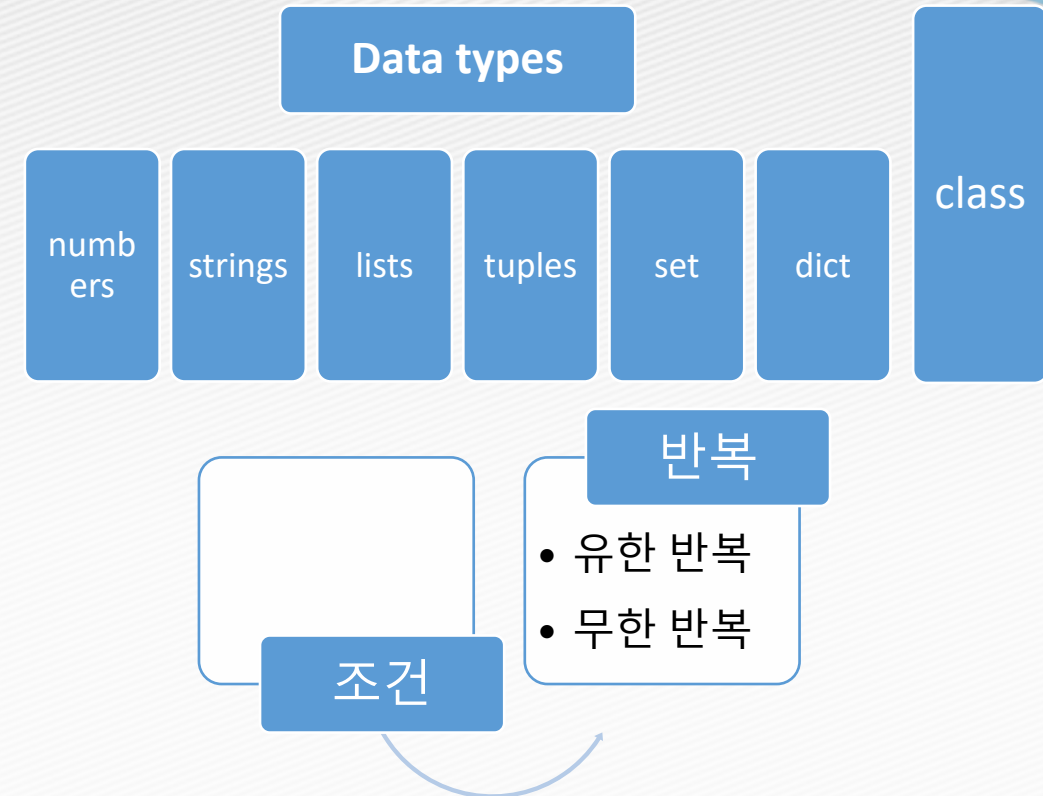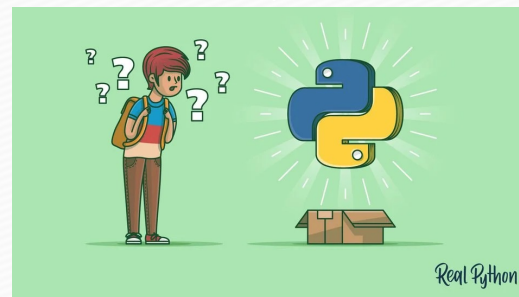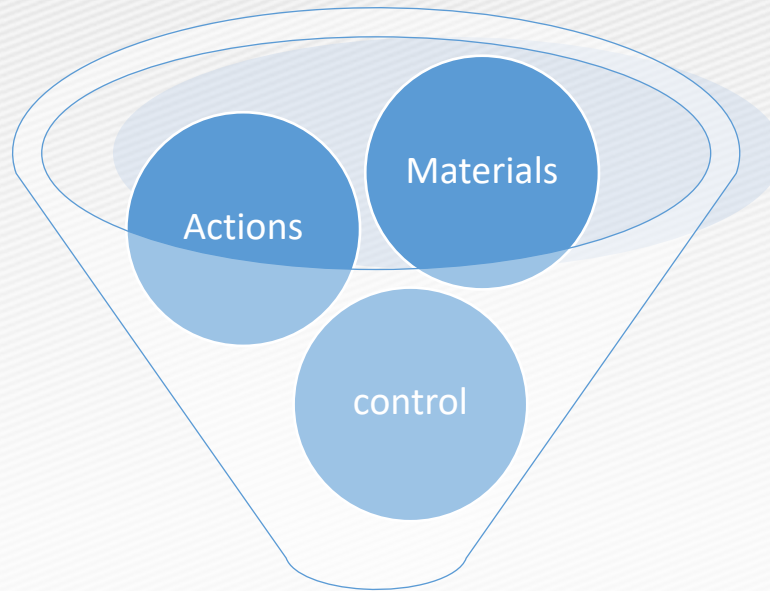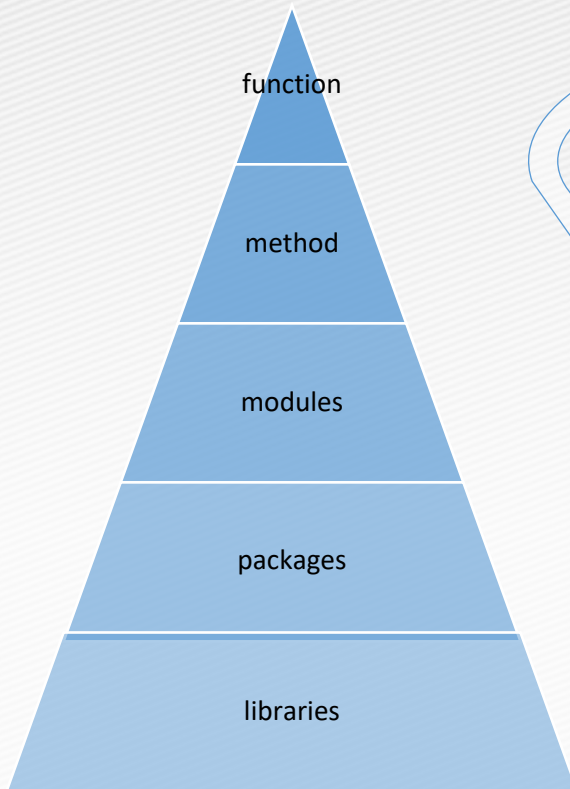1976–1987

Guido Van Rossum

**Guido van Rossum** ✔
@gvanrossum

Python's BDFL-emeritus, Distinguished Engineer at Microsoft, Computer History Fellow. Opinions are my own. He/him.

📍 San Francisco Bay Area   🔗 python.org/~guido/   📅 가입일: 2008년 8월

**515** 팔로우 중   **19.9만** 팔로워

팔로우

https://www.python.org/doc/essays/cp4e/

# 실습

## `pip3 install wordcloud`  <span style="background:red;color:white">library</span>

text = "Sungshin Sungshin Sungshin Sungshin University University"
<span style="background:red;color:white">variable</span>   <span style="background:red;color:white">str(ing) data type</span>

import wordcloud
<span style="background:red;color:white">library</span>

wc_obj = wordcloud.**WordCloud()**
<span style="background:red;color:white">instance</span>   <span style="background:red;color:white">class</span>

wc = wc_obj.generate(text)
<span style="background:red;color:white">method</span>

wc.words_
<span style="background:red;color:white">attribute</span>
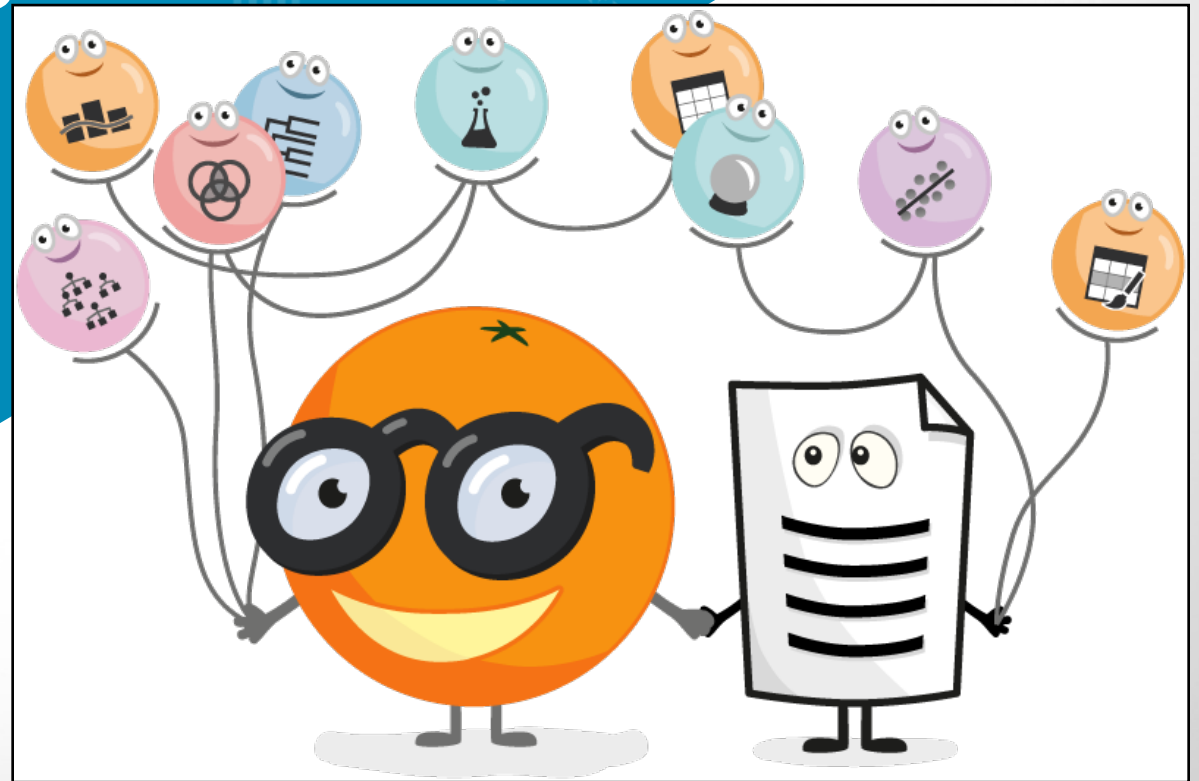
# Part 2
# Installing orange

Orange is developed by <u>Bioinformatics Lab</u> at University of Ljubljana, Slovenia.

**<u>Initial release</u>:** 10 October 1996; 26 years ago

# Orange 설치



or

# Orange 설치

Windows                    macOS

<span style="color:red">or</span>

conda config --add channels conda-forge          pip install PyQt5 PyQtWebEngine
conda install pyqt                                pip install orange3
conda install orange3

# Orange Visual Programming



https://orange3.readthedocs.io/projects/orange-visual-programming/en/latest/index.html

# Part 3
# Visual Programming with orange

# Workflow for a data analysis project



How domain expertise can help us during a data analysis project?

**Project Objective**
What questions should we ask?

**Data Preprocessing**
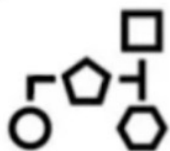How they relate to our project objective?

**Model Training**
What are the rules of thumb for sanity checks? How can we improve the training?

**Data Collection**
What variables do we need?
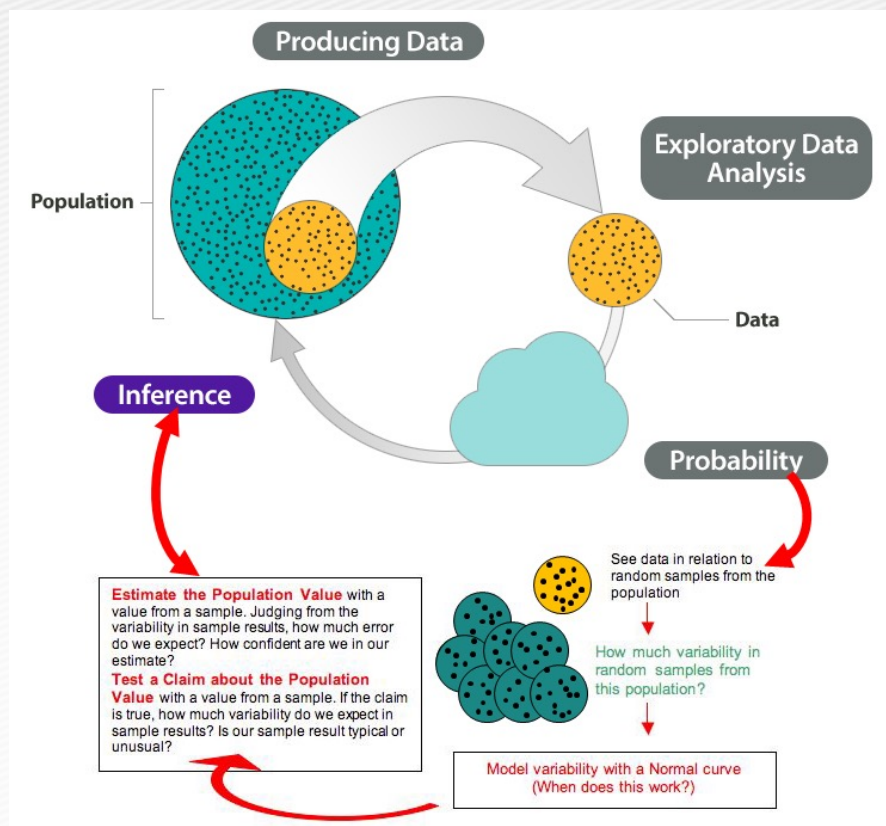
**Model Selection**
What would be the most appropriate model(s)?

**Results Interpretation**
What results are more useful or actionable?

Source: https://blog.ml.cmu.edu/2020/08/31/1-domain-knowledge/

# Statistical Inference vs. Modeling

## Inference

## Modeling



https://www.coursesidekick.com/statistics/study-guides/wmopen-concepts-statistics/wim-linking-probability-to-statistical-infe

# Classic Example: iris classification

https://archive.ics.uci.edu/ml/datasets/Iris

붓꽃의
품종



Measurement (in cm)

# Orange: Building workflow by connecting widgets



Data Table outputs selected data. Since no data is selected, the connection is empty (dashed). Scatter Plot is empty.

Scatter Plot show a plot. The connection is full, as the File widget is sending data directly to the widget.

# How to build in Orange

# Loading data using File widget



- Loading dataset we have to use **File Widget** which is available in Data Section.

- After clicking on File Widget, it will automatically appear on Canvas then you have to double click on that widget.
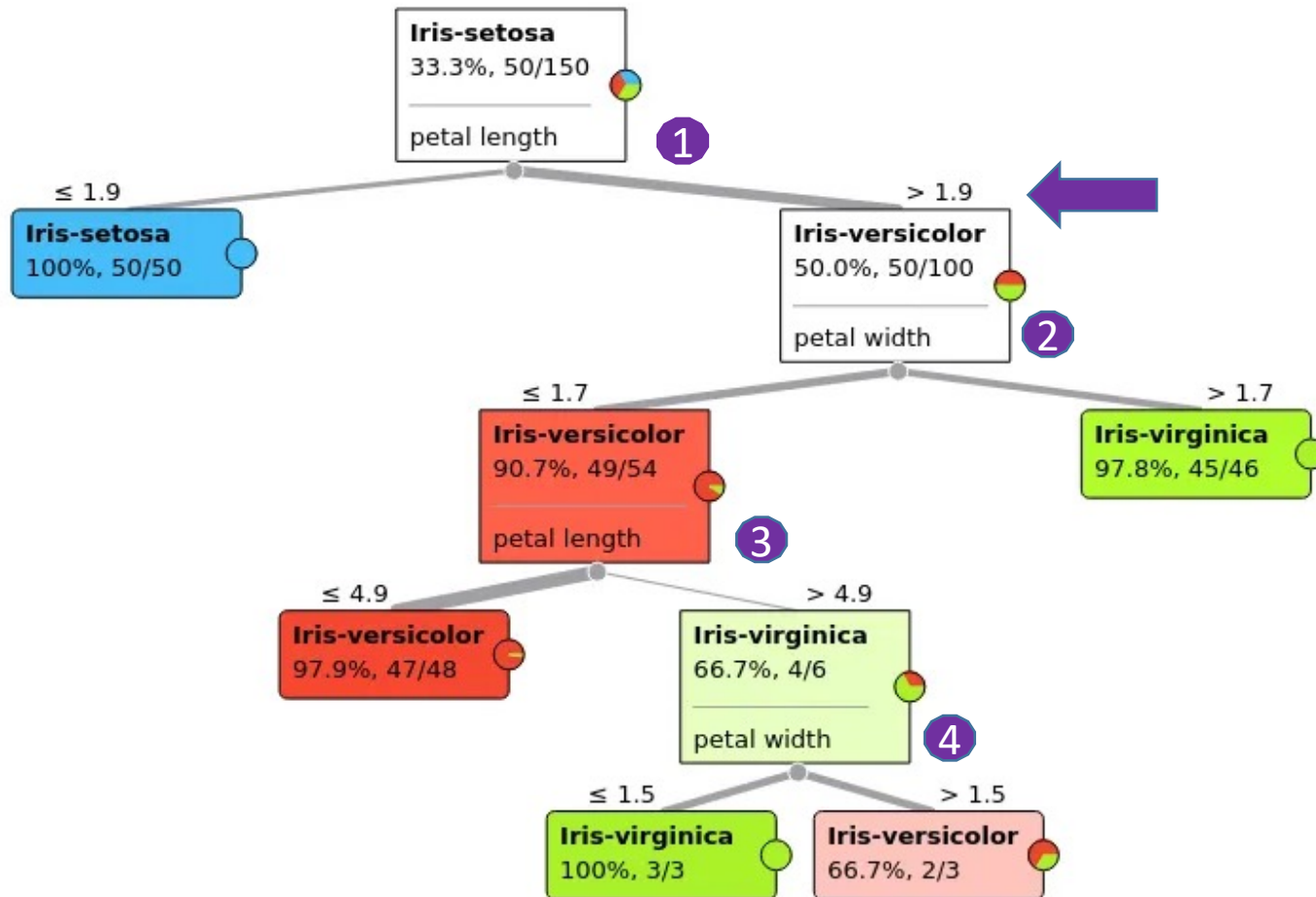
# Selection subset of the data (1)

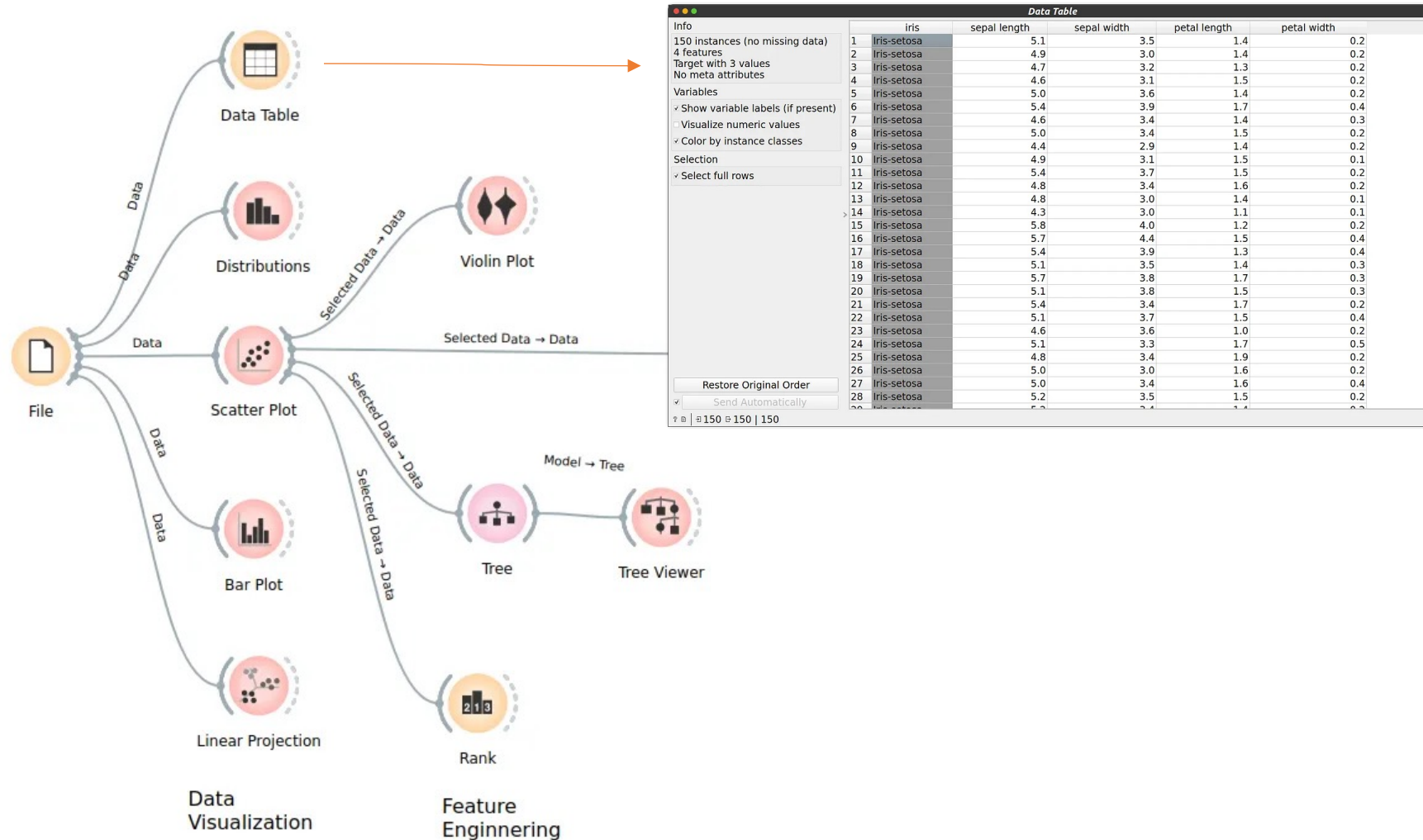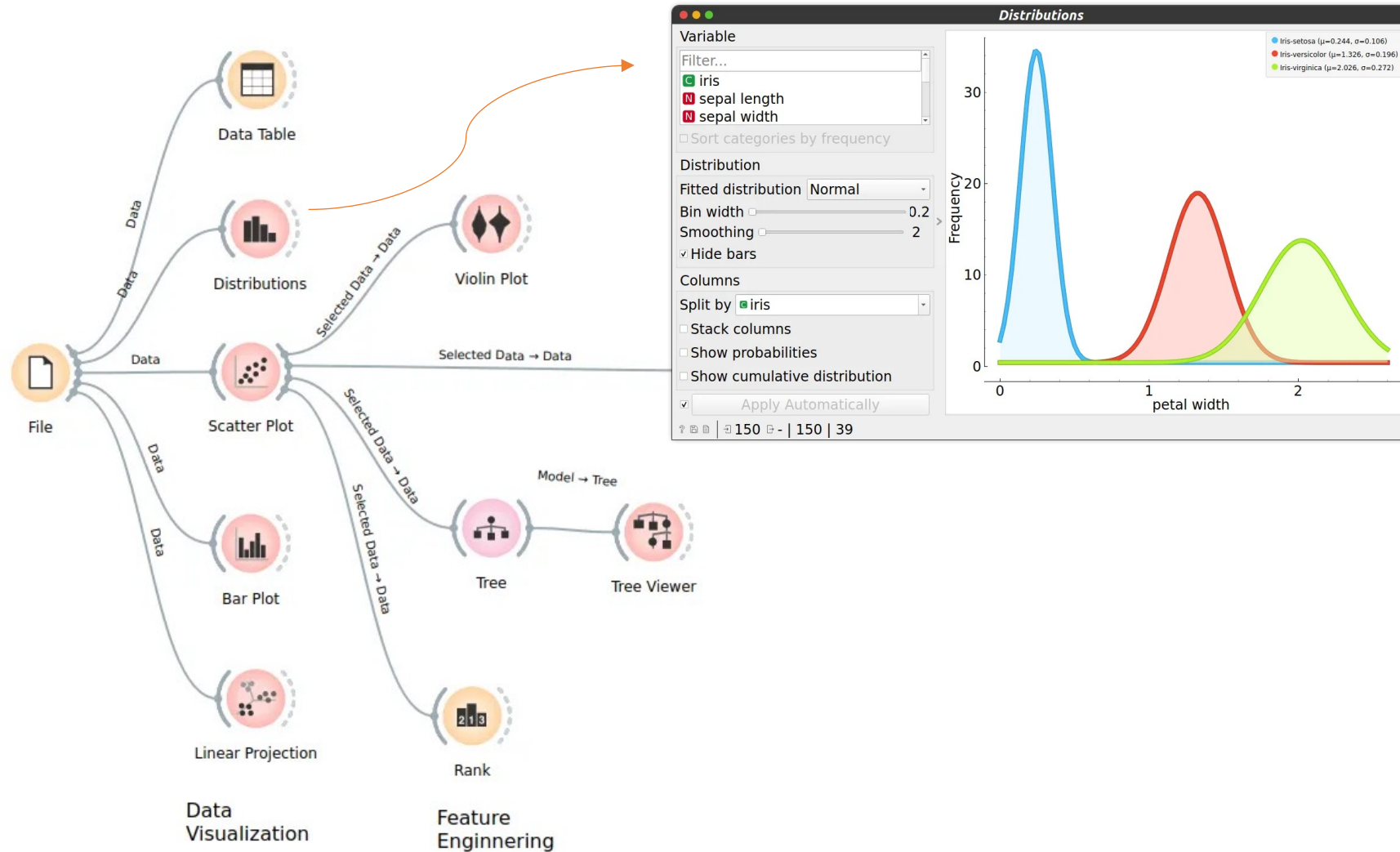# Selection subset of the data (2)

# Can we create this workflow example?
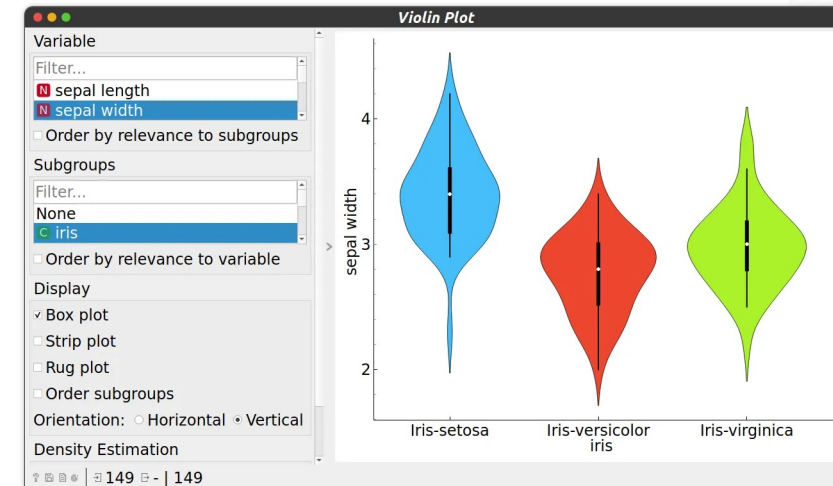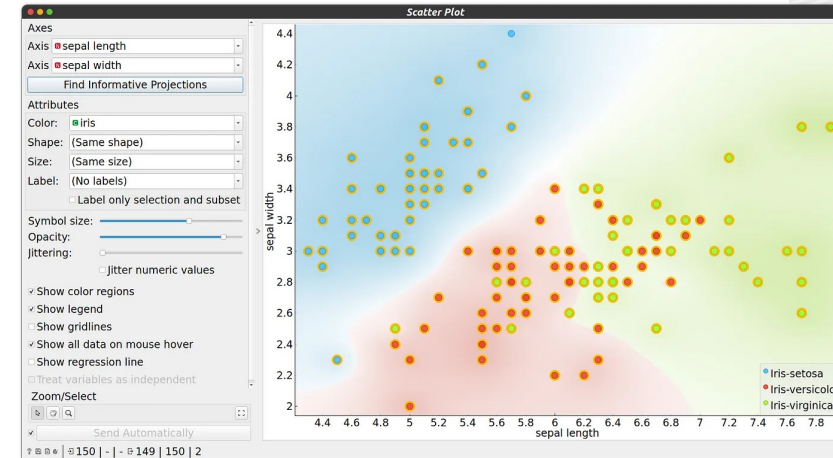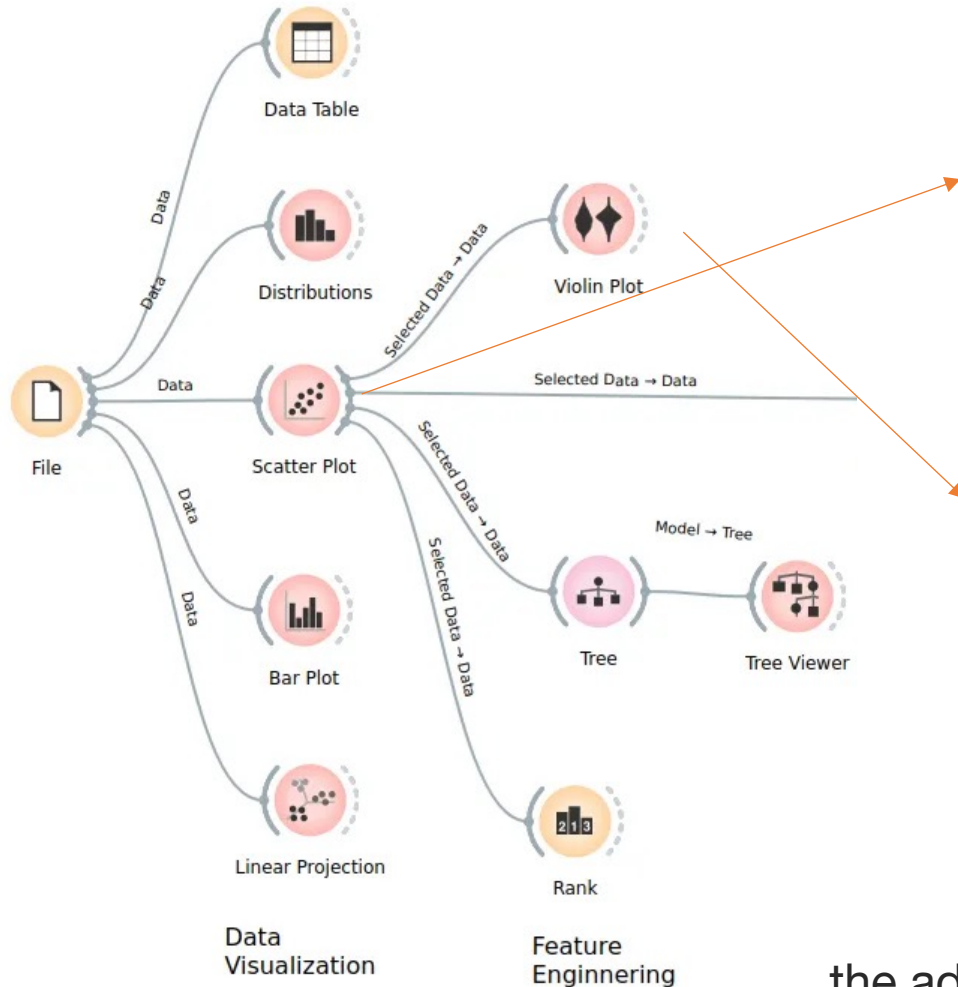
# How to read the output of classification tree

# Data Visualization – Data Table

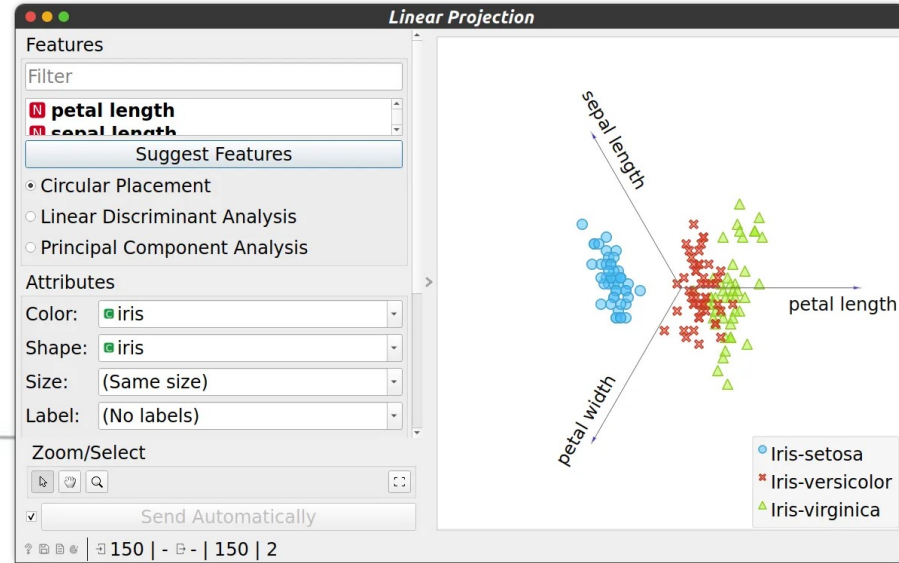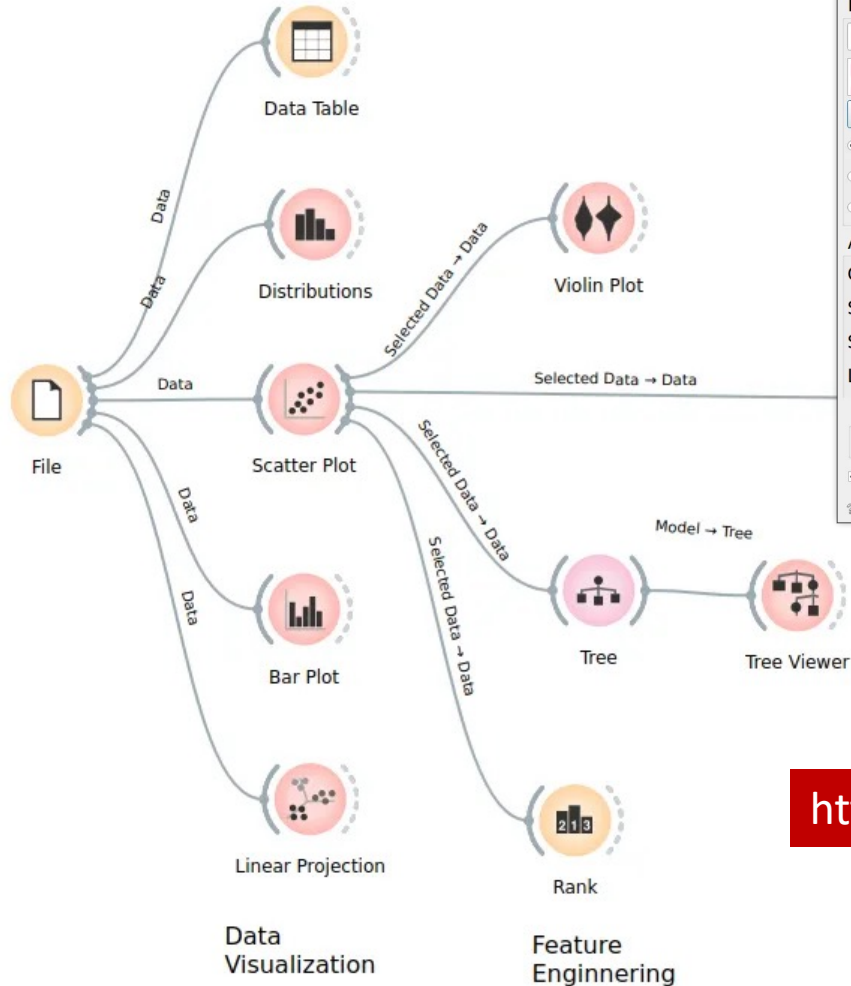# Data Visualization – Distribution

# Data Visualization – Scatter & Violin Plots



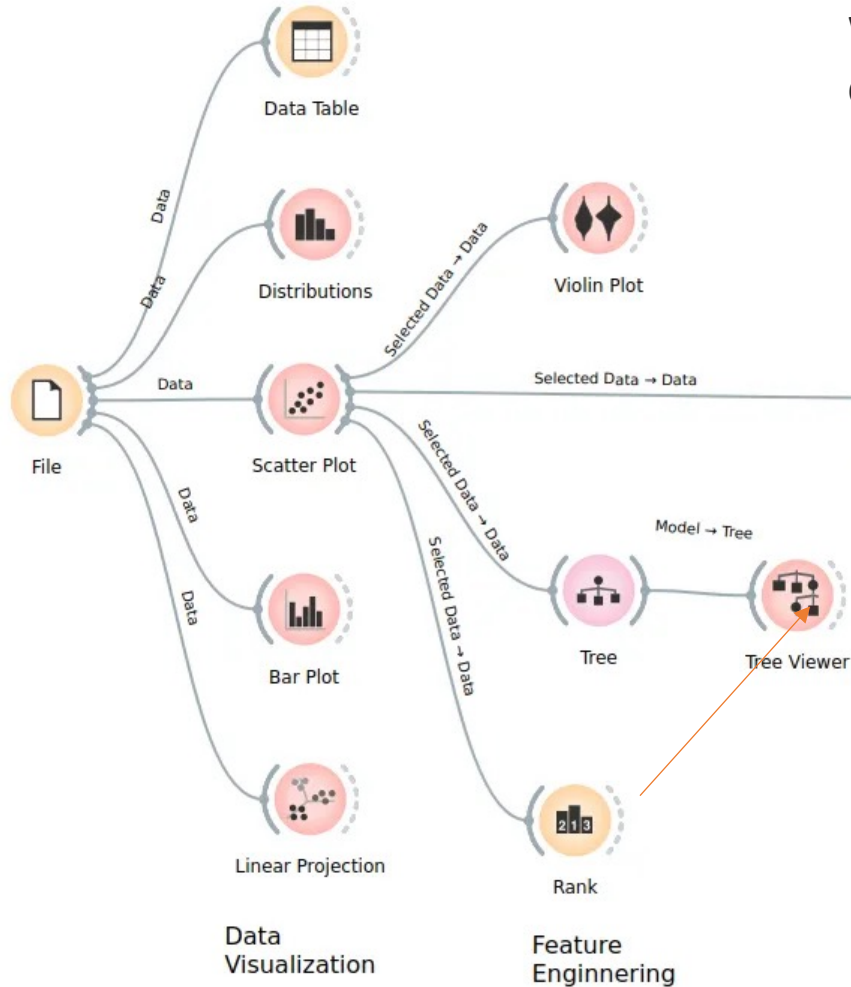the addition of a rotated **kernel density plot** on each side
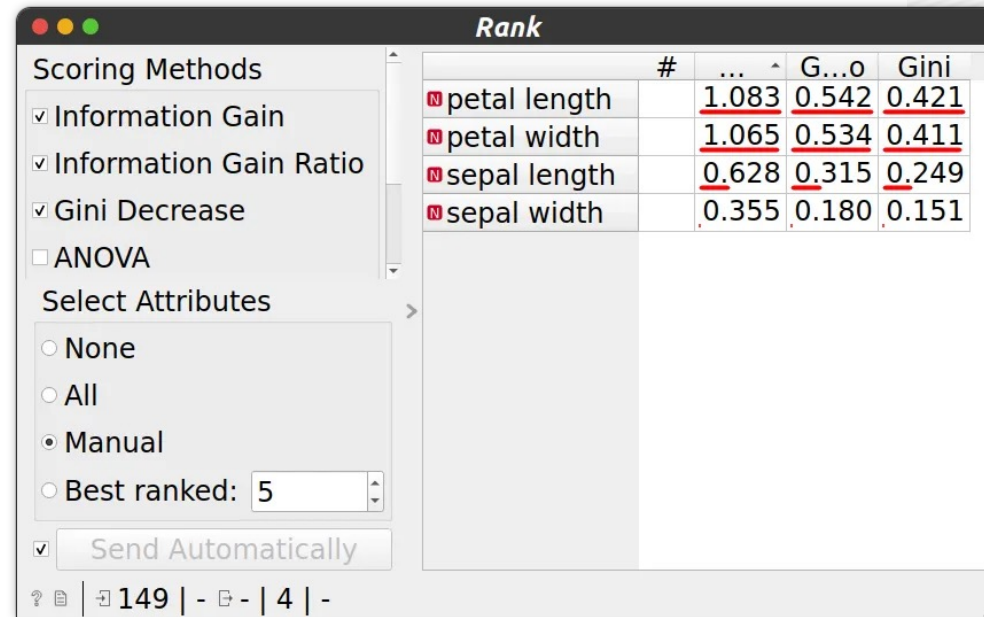
# Data Visualization – Linear Projection



visualize the data up to 3D

https://orangedatamining.com/widget-catalog/visualize/linearprojection/

# Data Visualization



Which features are most important for classification

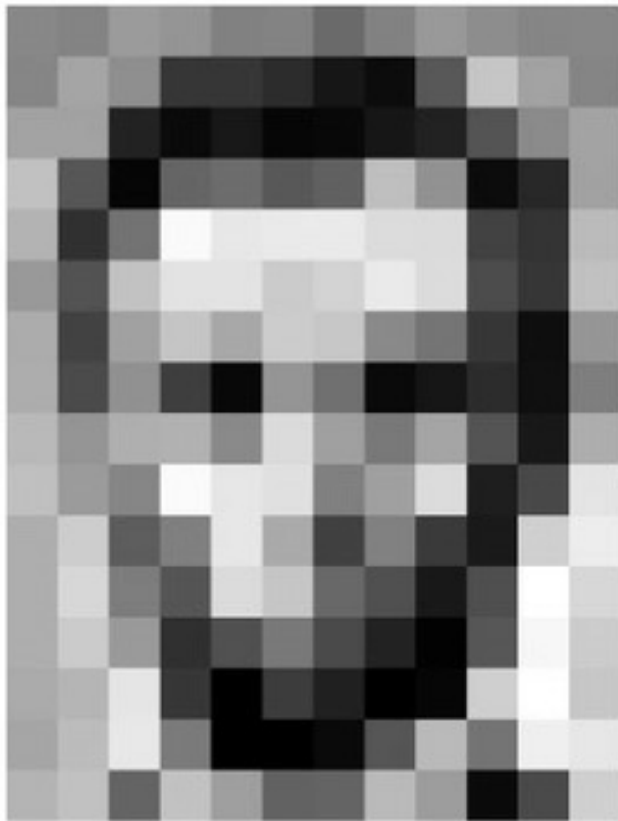**If you are interested in what widgets are available, click the blue circle above.**

# Part 4
# Word Cloud

2:00–2:20

# Image as numbers

- 흑백 그림은 픽셀의 밝기 값을 0-255 사이의 값으로 표현한 이미지

pixel



https://www.researchgate.net/figure/mage-of-Abraham-Lincoln-as-a-matrix-of-pixel-values_fig1_330902210
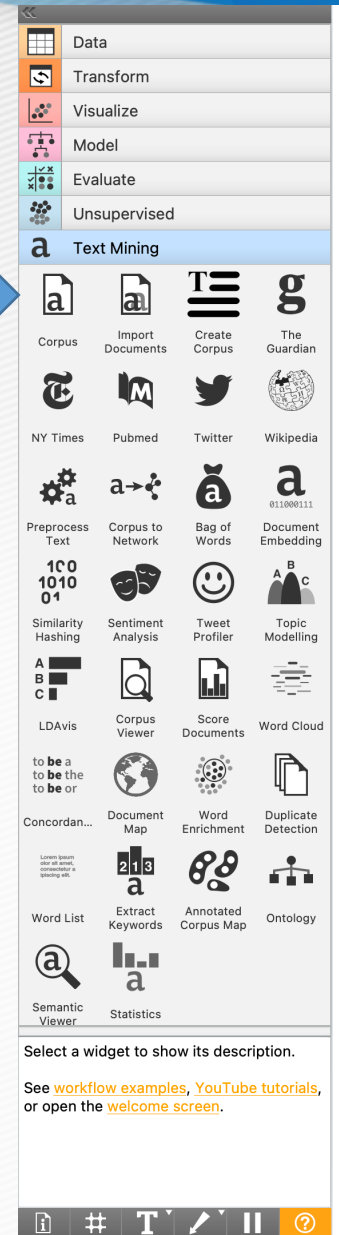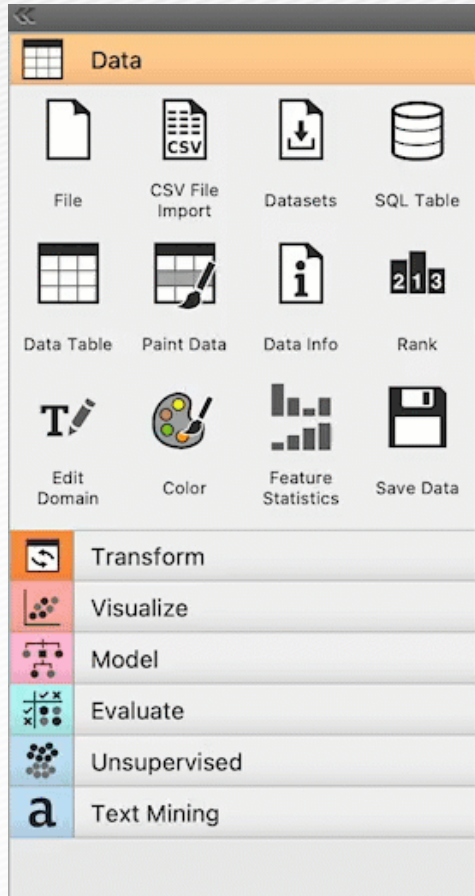
# Making Word Cloud using CLI

```
from wordcloud import WordCloud
text = "Sungshin Sungshin Sungshin Sungshin \
Sungshin Sungshin Sungshin Sungshin Sungshin Sungshin \
English English English English English Department"
wc = WordCloud().generate(text)
print(dir(wc))
```

```
import matplotlib.pyplot as plt
plt.figure()
plt.imshow(wc)
plt.axis("off")
plt.show()
#plt.savefig("wc_sungshin.png")
```
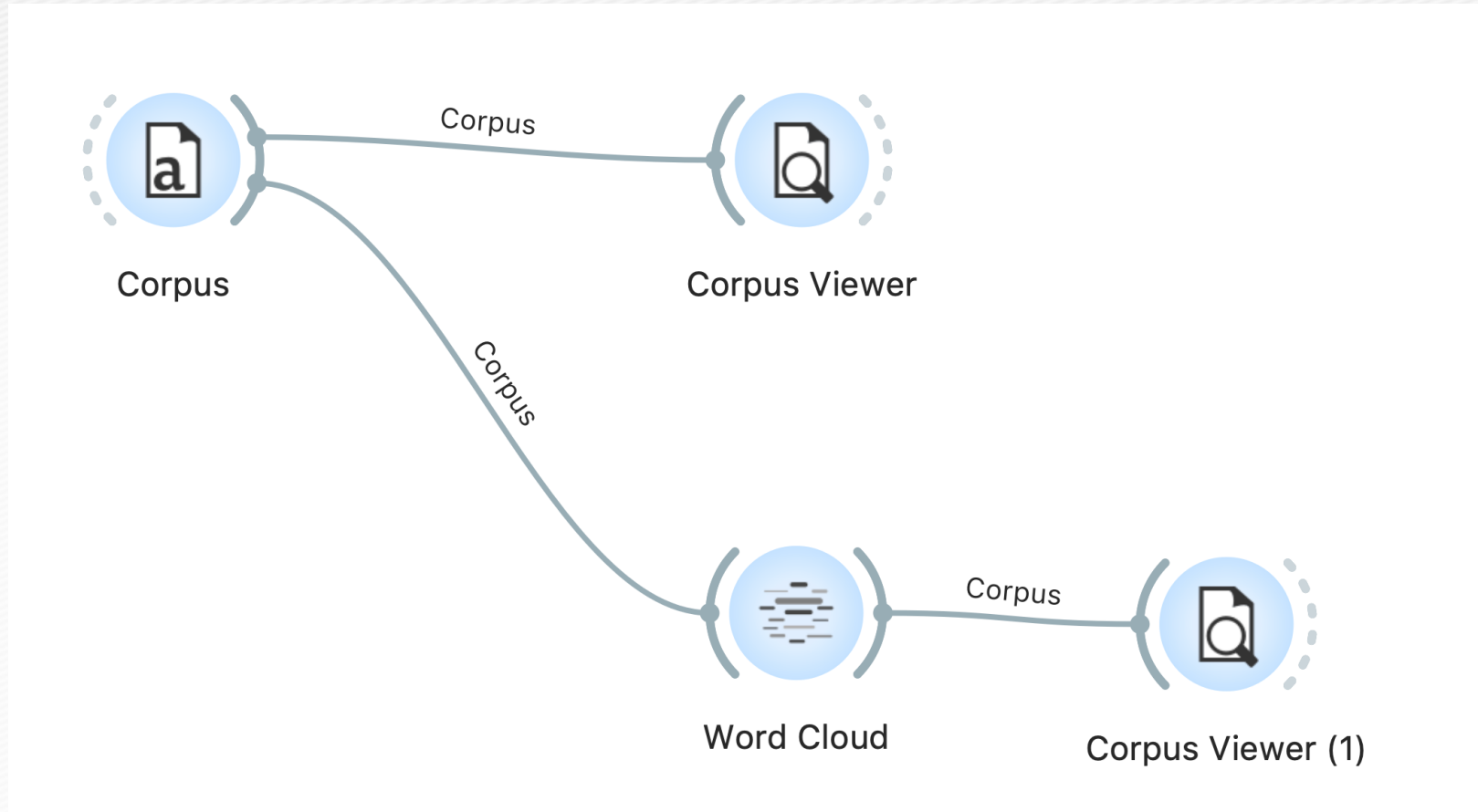
# Options > Add-ons… > Text



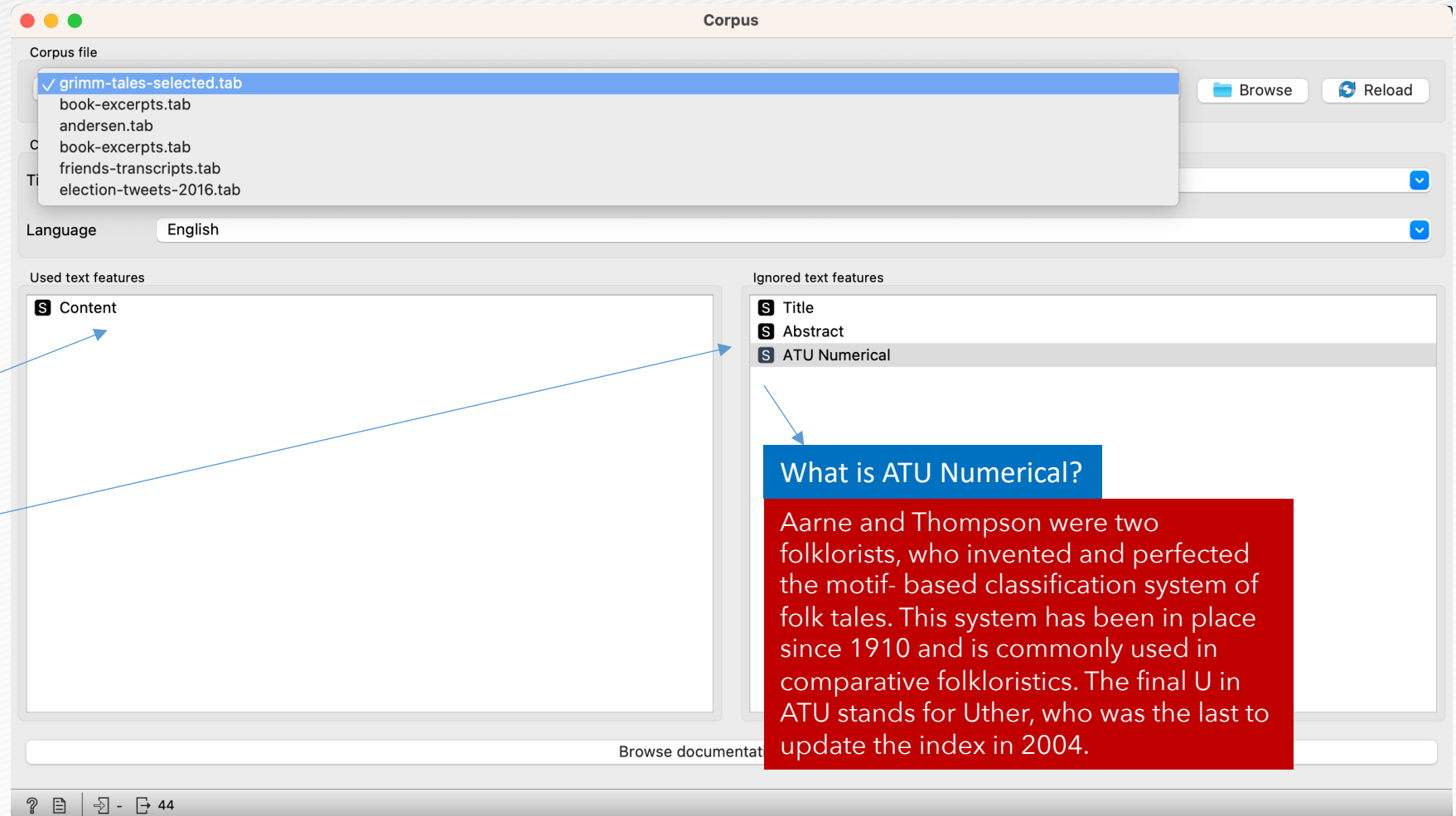Corpus is any collection of documents.

# Workflow for Word Cloud

- Start by constructing a workflow that consists of a Corpus widget, a Word Cloud widget and two Corpus Viewer widgets:
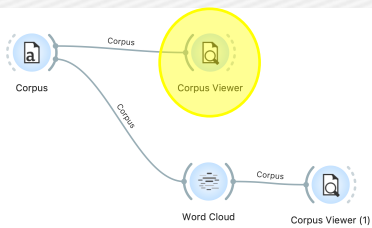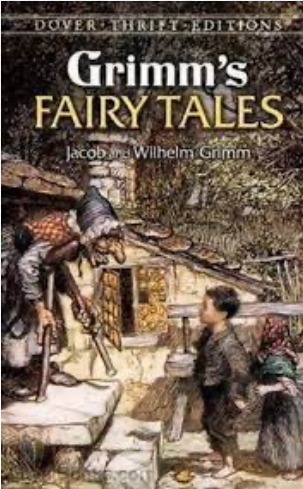
# Preloaded dataset

- From "Browse documentation data sets..." choose *Grimm-tales-selected.tab*, a data set containing Grimm's selected tales.

  - The particularity of the Corpus widget is that it sets the text feature(s) to apply text mining on.
  - "Used text features" defines the content (text), while other columns contain meta attributes (title, abstract, etc.).

# Data: Grimm-tales-selected (44 selected Grimm's tales)

- preloaded dataset



https://en.wikipedia.org/wiki/Grimms%27_Fairy_Tales

# Word Cloud via Orange

- Open Word Cloud.
- Word Cloud displays word frequencies, where the more frequent the word, the larger the font.

# Text cleaning needed!

- This word cloud is a mess! We got a bunch of semantic junk in our visualization. Is there a way to clean this up?



Preprocessing

# Part 5
# Text Preprocessing

# Text Preprocessing

- In the Preprocess Text widget, we decided to transform all words to lowercase, treat each word as a token (and omit punctuation), and to remove the stopwords (such as "in", "and", and "the").

- This preprocessing outputs the following tokens:

  "This is a sample sentence." → "sample", "sentence"

- To see the results of preprocessing, we can display the most frequent tokens in Word Cloud. Word Cloud enables us to identify redundant words and irregularities.

# Word Cloud after preprocessing

We see the results of our preprocessing in the Word Cloud. Two of the most frequent words are "would" and "could". If we decide these two words are not important for our analysis, it would be good to omit them.

# Custom filtering

- Load the list of custom stopwords in the right-hand dropdown of the Filtering section.



A good plain text editor is Sublime, but you can easily work with Notepad (or 메모장).

# Other ways of filtering

**Regular Expression**

https://www.w3schools.com/python/python_regex.asp

https://regexr.com

# Alice's Adventures in Wonderland

```python
from wordcloud import WordCloud
text = open('alice.txt').read()
wc_obj = WordCloud()
wc = wc_obj.generate(text)

import matplotlib.pyplot as plt

plt.figure()
plt.imshow(wc)
plt.axis("off")
plt.show()
plt.savefig("wc_alice.png")
```

# STOPWORDS in the wordcloud library

```python
from wordcloud import WordCloud, STOPWORDS
stopwords = set(STOPWORDS)
stopwords.add("said")
wc = WordCloud(stopwords=stopwords)
wc = wc.generate(text)
plt.figure(figsize=(12,12))
plt.imshow(wc)
plt.axis("off")
plt.show()
```

# Masked Word Cloud만들기



준비물: alice_mask.png

```
from PIL import Image
import numpy as np
alice_mask = np.array(Image.open('alice_mask.png'))
plt.figure()
plt.imshow(alice_mask)
plt.axis("off")
```

```
array([[255, 255, 255, ..., 255, 255, 255],
       [255, 255, 255, ..., 255, 255, 255],
       [255, 255, 255, ..., 255, 255, 255],
       ...,
       [255, 255, 255, ..., 255, 255, 255],
       [255, 255, 255, ..., 255, 255, 255],
       [255, 255, 255, ..., 255, 255, 255]], dtype=uint8)
```

```
wc_obj = WordCloud(mask=alice_mask)
wc = wc_obj.generate(text)

plt.figure()
plt.imshow(wc)
plt.imshow(alice_mask)
plt.axis('off')
plt.show()
```

# Part 7
# Loading your own data

2:35–3:00

# iris.csv

```
"sepal.length","sepal.width","petal.length","petal.width","variety"
5.1,3.5,1.4,.2,"Setosa"
4.9,3,1.4,.2,"Setosa"
4.7,3.2,1.3,.2,"Setosa"
```



**CSV File Import**

File: iris.csv

Info

150 rows, 5 features, 0 metas

Import Options... | Cancel | Reload

? 150

Encoding: Unicode (UTF-8)

Cell delimiter: Comma ,

Quote character: "

Number separators: Grouping: Decimal: .

Column type:

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | sepal.length | sepal.width | petal.length | petal.width | variety |
| 2 | 5.1 | 3.5 | 1.4 | .2 | Setosa |
| 3 | 4.9 | 3 | 1.4 | .2 | Setosa |
| 4 | 4.7 | 3.2 | 1.3 | .2 | Setosa |
| 5 | 4.6 | 3.1 | 1.5 | .2 | Setosa |
| 6 | 5 | 3.6 | 1.4 | .2 | Setosa |
| 7 | 5.4 | 3.9 | 1.7 | .4 | Setosa |
| 8 | 4.6 | 3.4 | 1.4 | .3 | Setosa |
| 9 | 5 | 3.4 | 1.5 | .2 | Setosa |
| 10 | 4.4 | 2.9 | 1.4 | .2 | Setosa |
| 11 | 4.9 | 3.1 | 1.5 | .1 | Setosa |
| 12 | 5.4 | 3.7 | 1.5 | .2 | Setosa |
| 13 | 4.8 | 3.4 | 1.6 | .2 | Setosa |
| 14 | 4.8 | 3 | 1.4 | .1 | Setosa |
| 15 | 4.3 | 3 | 1.1 | .1 | Setosa |
| 16 | 5.8 | 4 | 1.2 | .2 | Setosa |
| 17 | 5.7 | 4.4 | 1.5 | .4 | Setosa |
| 18 | 5.4 | 3.9 | 1.3 | .4 | Setosa |
| 19 | 5.1 | 3.5 | 1.4 | .3 | Setosa |

Reset | Restore Defaults | Cancel | OK

## Attribute Type

- C: Continuous
- D: Discrete
- T: Time
- S: String

# Loading a csv file



Statistical information like median, mode, Missing values, Distribution, Min, Max Values

For Loading the CSV from the local device we have to use **CSV Import** Widget

Features

Target

# Import your own text documents

# Part 8
# Concordance

3:00–3:10

# Visualizing corpus - Concordance

- We have already seen some of the preprocessing results in a word cloud.
  - Word Cloud shows us word frequencies.
- But we still don't know much about the use of a specific word in a text.
  - For example, 'oh' could be a lowercase version of OH (the chemical compound of hydroxide), a simple exclamation 'Oh!' or an abbreviation for the state of Ohio.
- To check the context of a particular word we can use Concordance widget.
  - Concordance shows us the text around our word.

# Concordance widget

- Connect Concordance to Corpus to pass the text to the widget.

# Concordance



- To browse the word, type it in the query line at the top or provide it with the Word Cloud.
- Here we have selected the word 'think' and observed the context in Concordance.

# Part 9
# Bag of Words

# Finding patterns in the text

- To find any patterns in our text, we need to convert documents into numeric vectors is to count the words in each text.

- Bag of Words creates a table with words in columns and documents in rows.

| | this | is | an | example | another | apple |
|---|---|---|---|---|---|---|
| "This is an example" | 1 | 1 | 1 | 1 | 0 | 0 |
| "Another example" | 0 | 0 | 0 | 1 | 1 | 0 |
| "This is another apple. | 1 | 1 | 0 | 0 | 1 | 1 |

# TF-IDF



$$TF(t,d) = \frac{\text{number of times t appears in d}}{\text{total number of words in d}} \qquad IDF(t) = \log \frac{\text{Total number of documents}}{\text{Number of documents that contain t}}$$

https://betterprogramming.pub/a-friendly-guide-to-nlp-tf-idf-with-python-example-5fcb26286a33

# TF-IDF toy example

**Review 1**: Game of Thrones is an amazing TV series!
**Review 2**: Game of Thrones is the best TV series!
**Review 3**: Game of Thrones is so great.

| Word | TF REV1 | TF REV2 | TF REV3 | IDF | TF-IDF REV1 | TF-IDF REV2 | TF-IDF REV3 |
|---|---|---|---|---|---|---|---|
| amazing | 1/8 | 0 | 0 | $\ln(4/2)+1=1.7$ | 0.47 | 0 | 0 |
| an | 1/8 | 0 | 0 | $\ln(4/2)+1=1.7$ | 0.47 | 0 | 0 |
| best | 0 | 1/8 | 0 | $\ln(4/2)+1=1.7$ | 0 | 0.47 | 0 |
| game | 1/8 | 1/8 | 1/6 | $\ln(4/4)+1=1$ | 0.28 | 0.28 | 0.32 |
| great | 0 | 0 | 1/6 | $\ln(4/2)+1=1.7$ | 0 | 0 | 0.54 |
| is | 1/8 | 1/8 | 0 | $\ln(4/4)+1=1$ | 0.28 | 0.28 | 0.32 |
| of | 1/8 | 1/8 | 1/6 | $\ln(4/4)+1=1$ | 0.28 | 0.28 | 0.32 |
| series | 1/8 | 1/8 | 0 | $\ln(4/3)+1=1.29$ | 0.36 | 0.36 | 0 |
| so | 0 | 0 | 1/6 | $\ln(4/2)+1=1.7$ | 0 | 0 | 0.54 |
| the | 0 | 1/8 | 0 | $\ln(4/2)+1=1.7$ | 0 | 0.47 | 0 |
| thrones | 1/8 | 1/8 | 1/6 | $\ln(4/4)+1=1$ | 0.28 | 0.28 | 0.32 |
| tv | 0 | 1/8 | 0 | $\ln(4/3)+1=1.29$ | 0.36 | 0.36 | 0 |

- the most present words, such as "game", "of", "thrones", "is", have the smallest IDF
- terms like "amazing" and "great" have higher TF-IDF values
- common words, like "so" and "the", contribute more too since they aren't present in all the sentences and we didn't remove the stop words to keep the approach as simple as possible.

# A Bag of Words Widget

Pass the data through **a Bag of Words widget** and then again to a Data Table.

We get a new column that contains word counts for each document.



Data: andersen.tab

# Part 10
# Clustering and Distances

# Distance and similar documents

- One common task in text mining is finding interesting groups of similar documents.

- That is, we would like to identify documents that are similar to each other. .

- We pass the data to Distances, use Euclidean distance, then to Hierarchical Clustering.

# Types of Distance



The Euclidean Distance is the shortest distance between two points.

The Manhattan distance is the sum of the lengths of the rectangle formed by the two points

the Cosine distance is the angle subtended at the origin between the two documents.

https://nickgrattan.wordpress.com/2014/06/10/euclidean-manhattan-and-cosine-distance-measures-in-c/

# Jaccard Similarity coefficient or Jaccard Index

- For text, an intuitive approach for measuring similarity would also be the number of words that two documents share.



$$J(A, B) = \frac{|A \bigcap B|}{|A \bigcup B|}$$

- The measure is called Jaccard similarity coefficient or Jaccard index.

- Note that in this case, we are measuring similarity, not distance.

# Workflow

- Now, let us go back to our Grimm's Tales and construct the following workflow:



**You can try the same workflow on a different corpus, say *bookexcerpt.tab*, which contains excerpts from adult and children's books.**

# Hierarchical Clustering

- **Ward's method**
  - *Least increase in total variance (around cluster centroids)*

- **Average linkage**
  - *Average distance between clusters*

- **Complete linkage**
  - *Max distance between clusters*

https://python-data-science.readthedocs.io/en/latest/_images/aggocluster2.png

# Hierarchical Clustering



Why are some Tales of Magic mixed with Animal Tales? What do they have in common?

# Part 11
# Classification

# predict the ATU type based on the content of the tale

- The Aarne-Thompson type (ATU)
  - the index of folk-tale motifs
  - It is already marked every tale with a high-level (genre) and a mid-level ATU type (subgenre).
- Could we perhaps predict the ATU type based on the content of the tale?



Aarne and Thompson were two folklorists, who invented and perfected the motif- based classification system of folk tales. This system has been in place since 1910 and is commonly used in comparative folkloristics. The final U in ATU stands for Uther, who was the last to update the index in 2004.

# Logistic Regression

$$p(X) = \beta_0 + \beta_1 X.$$

$$\log\left(\frac{p(X)}{1-p(X)}\right) = \beta_0 + \beta_1 X.$$

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}.$$



Linear Regression



Logistic Regression

https://medium.com/data-science-group-iitr/logistic-regression-simplified-9b4efe801389

# Prediction

- Target – an ATU type

- Feature – Numerical representation of each document

- Model – Logistic Regression

- Prediction - a column with predicted values from Logistic Regression

# Part 12
# Predictions with test
# data

# Data splitting



Train,Test Split Dataset

Serve the same Purpose for K-Fold,Cross Validation

| Training Data | Holdout | Holdout |
|---|---|---|
|  | Validation Data (Tuning Hyper-parameter) | Testing Data (Evaluating Performance) |

| Training Data (for Fitting) | Holdout |
|---|---|
|  | Testing Data (Evaluating Performance) |

Original Available Data

Jesus Saves @JCharisTech

# Predicting on new data

- Predicting on new data works just like for regular data.

# Guess game – guess the tale type

- Open a new Corpus widget and load the ***andersen.tab*** corpus.
  - Three tales from H. C. Andersen.

- Connect them to Predictions the same way as before - with Logistic Regression passing the constructed model and the new Corpus widget passing the data for prediction.

- Logistic Regression predicted two tales to be Tales of Magic and one the Animal Tale.

| | Logistic Regression | Title | Content |
|---|---|---|---|
| 1 | 0.01 : 0.99 → Tales of Magic | The Little Match-Seller | It was terribly cold and nearly dark on... |
| 2 | 0.00 : 1.00 → Tales of Magic | The Philosopher's Stone | Far away towards the east, in India, w... |
| 3 | 0.90 : 0.10 → Animal Tales | The Ugly Duckling | It was lovely summer weather in the c... |

# Resources

4:30–5:00

# References

- Orange Lecture Notes
  - https://orangedatamining.com/blog/2020/2020-02-08-lecture-notes/
- New Video Tutorials on Text Mining
  - https://orangedatamining.com/blog/2020/2020-09-28-text-tutorials/
- Observing Word Distribution
  - https://orangedatamining.com/blog/2021/2021-01-27-word-distribution/
- Machine Learning Jargon
  - https://orangedatamining.com/blog/2022/2022-02-01-machine-learning-jargon/
- Semantic Analysis of Documents
  - https://orangedatamining.com/blog/2021/2021-09-17-semantic-analysis/
- PCA vs. MDS vs. t-SNE
  - https://orangedatamining.com/blog/2021/2021-06-17-pca-mds-tsne/
- How to identify fake news with document embeddings
  - https://orangedatamining.com/blog/2020/2020-10-15-document-embedders/
- Detecting Story Arcs with Orange
  - https://orangedatamining.com/blog/2020/2020-07-27-story-arcs/