

Data Science Infrastructures – Exercise 03 (DSI E03 ST 24)

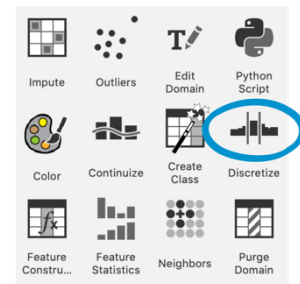
02/05/2024 / Philipp Wieder / philipp.wieder@gwdg.de

General Prerequisites

To fulfil this assignment, you need the Software Orange, which you can download from <https://orangedatamining.com/>. Please follow the installation instructions for your respective operation system (<https://orangedatamining.com/download/>). Orange is freely available and comes with a number of so-called widgets that offer user interfaces for particular data science tasks. It has a comprehensive documentation including tutorials (<https://orangedatamining.com/docs/>) and comes with a number of example workflows (<https://orangedatamining.com/examples/>).

Assignment I – Data Type Portability: Discretization

As part of the Pre-Processing lecture we came across a step called data type portability, which is used to convert between different data types to make it easier to process certain data sets or apply standard algorithms. The conversion of numeric data types into categorical ones is called discretization. Orange has a particular widget for this task.



Select a Data Set

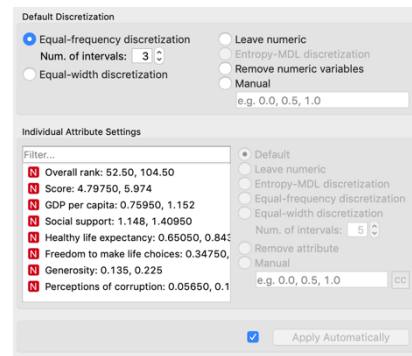
| Name | Type | Role | Values |
|--------------------|-----------|---------|--------|
| 1 Overall rank | N numeric | feature | |
| 2 Score | N numeric | feature | |
| 3 GDP per ... | N numeric | feature | |
| 4 Social support | N numeric | feature | |
| 5 Healthy life ... | N numeric | feature | |
| 6 Freedom to ... | N numeric | feature | |
| 7 Generosity | N numeric | feature | |
| 8 Perceptions ... | N numeric | feature | |
| 9 Country or ... | S text | meta | |

The first thing to be done is to select a data set. In the exercise folder, you can find the "World Happiness Report 2019.csv"¹. Select the "File" widget, choose the aforementioned file from wherever you have stored it on your machine, and it will be loaded into your widget. Double-clicking on the icon in Orange opens a dialogue that provides basic information about the data (like e.g. the number of instances and features) and lists all features and their respective data types. In our example, all data is numeric except the name of the country.

¹ Downloaded from <https://www.kaggle.com/datasets/unsdsn/world-happiness?select=2019.csv>.

Apply Discretization

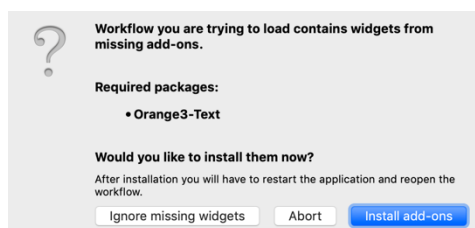
Now choose the “Discretization” widget and link it to the “File” widget. You can either do this by clicking on the dotted curve at the right of the “File” widget and then select the “Discretization” widget from the drop-down list or you place the “Discretization” widget directly on the main panel of Orange and then link both widgets. Just make sure that they are in the right order, i.e. the output from the “File” widget (which is the data set) is the input to the “Discretization” widget. Per default, equal-frequency discretization with three intervals is chosen. This splits each feature/attribute into three intervals, so that they each contain approximately the same number of instances. You can change the number of intervals or chose another discretization method, either for all attributes or individually². It does not necessarily make sense to discretize all attributes, so you can also leave it numeric.



Assignment

Discretize the attribute “Overall rank” into the Top 10, Top 11-50, and Top 51-100 of the happiest countries in the world in 2019. Leave all other attributes as is (you can apply settings to a selection of attributes at once). Look at the results and compare it to the original data set (you can use the “Data Table” widget for that). Then use the “Select Rows” widget to select the Top 10 and look at the results. Make screenshots of the complete workflow, the discretization settings, and the resulting data.

Assignment II – Pre-Processing of Text Data



For this assignment, you need a widget collection called “Text Mining”. One way is to download a text mining workflow example (e.g. “620-text-clustering.ows” from the exercise folder on Stud.IP³), because Orange then recognizes the missing widgets and asks whether you want to install them. Install the add-ons and you have everything you need for the assignment. But you can also install it using the

“Add-ons” manager from the “Options” menu.

² More info can be found here: <https://orange3.readthedocs.io/projects/orange-visual-programming/en/latest/widgets/data/discretize.html>.

³ More Workflows can be found here: <https://orangedatamining.com/examples/>.

If you are interested, try out the other pre-processing steps or use different settings.