

Supervised Learning Classification of yeast dataset

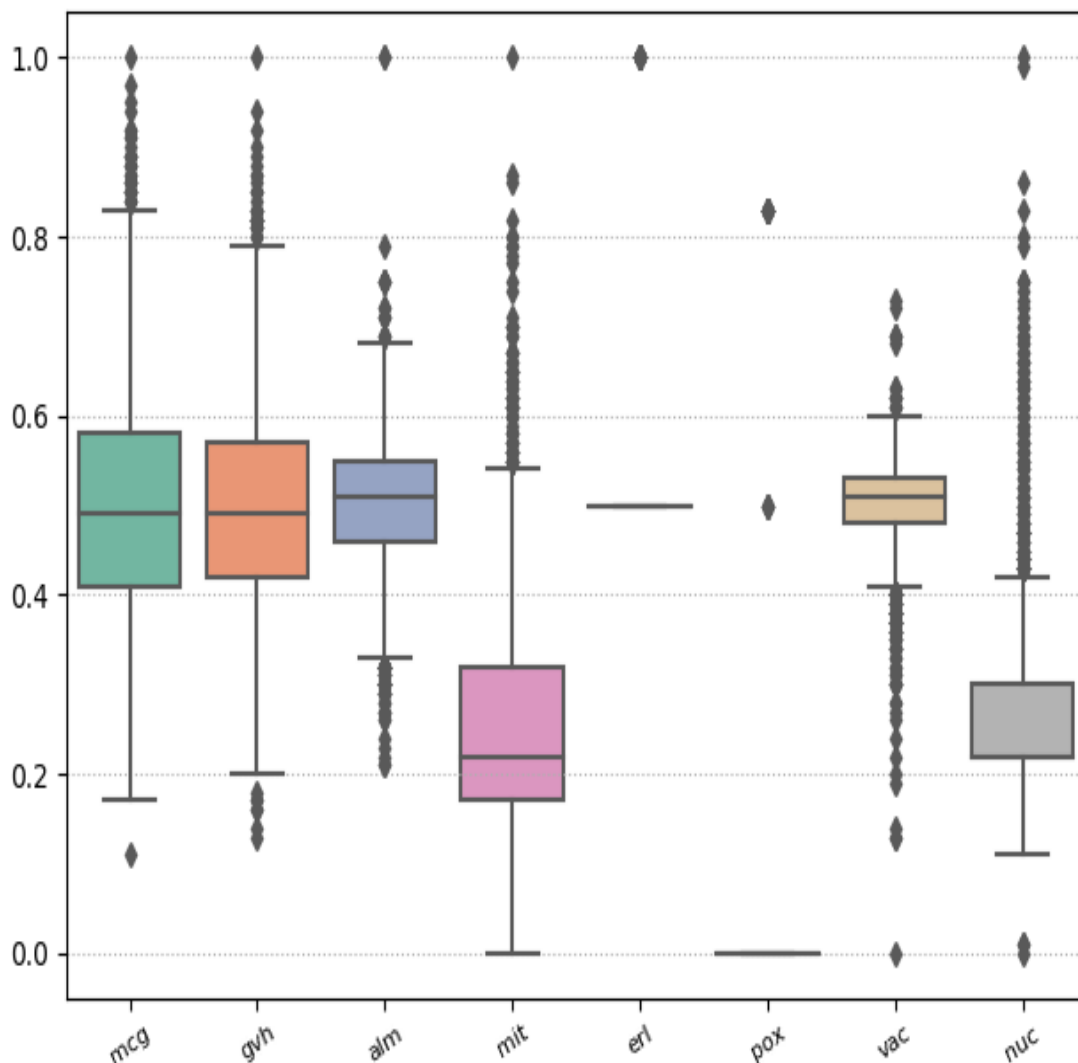
Content:

- 1) Get an overview of dataset
- 2) Compare Classification algorithms

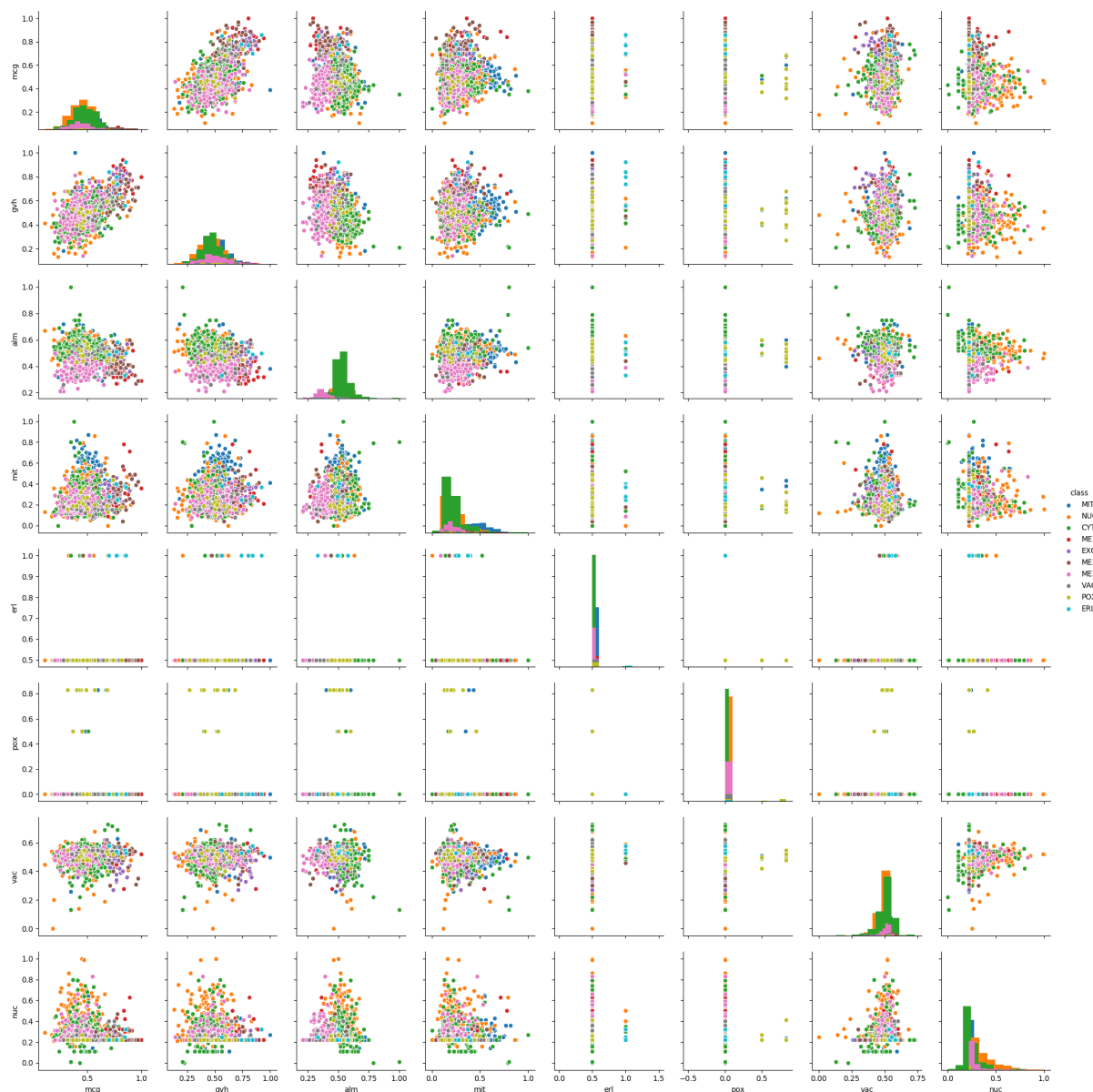
The following report shows a comparison of different supervised learning algorithms that classify the data of the yeast dataset.

1) *Get an overview of the dataset:*

The boxplot diagram shows an overview of the different features within the dataset and its attributes. On the X-Axis the different features of the yeast dataset are represented. The Y-Axis shows how the attributes of the different attributes of the features. The boxplot is based on the minimum, first quartile, median, third quartile, and maximum.



The second graph is a scatterplot that shows how the feature attributes are distributed for the different classes. Classification algorithms should show better results if clusters can be recognized within the scatterplot.



Now that we got a feeling of how the dataset is distributed, a random split for training and evaluation data has been conducted. The following boxplot diagram shows the comparison of the different supervised learning algorithms:



LDA: 0.553043 (0.049513)

CART: 0.468251 (0.050018)

SVM: 0.381441 (0.065492)

Figure 1 displays confusion matrices for 10-fold cross-validation of 10 models. The models are arranged in two rows: LR (k=10), LDA (k=10), KNN (k=10) in the top row, and CART (k=10), NB (k=10), SVM (k=10) in the bottom row. The y-axis represents 'predictions' and the x-axis represents 'validations', both with categories: CYT, ERL, EXC, ME1, ME2, ME3, MIT, NUC, POX, and VAC. The color scale ranges from 0 (dark green) to 10 (red).