

# Data preprocessing

Marcin Kuta

# Data imputation

- removal of columns or rows with missing values
- imputation of missing values with mean, median or mode
- imputation of missing values with the most frequent values, zero value or random value
- imputation of missing values with k-NN method

# Feature encoding

- one-hot encoding

# Feature scaling

- normalization (max-min scaling)

$$x \leftarrow \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

- standardization

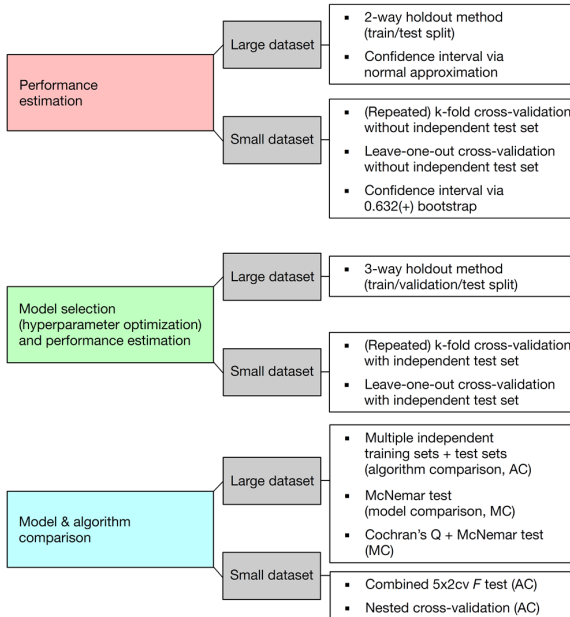
$$x \leftarrow \frac{x - \mu}{\sigma}$$

- soft-max scaling

$$y \leftarrow \frac{x - \mu}{r\sigma}$$

$$x \leftarrow \frac{1}{1 + \exp(-y)}$$

# Evaluation



- [1] <https://colab.research.google.com/github/jakevdp/PythonDataScienceHandbook/blob/master/notebooks/05.04-Feature-Engineering.ipynb>
- [2] <https://github.com/rasbt/machine-learning-book/tree/main/ch04>
- [3] Sebastian Raschka,  
Model Evaluation, Model Selection and Algorithm Selection in  
Machine Learning, 2018.