# Machine Learning in Python
## Supervised Learning - Classification and Metrics

Cristian A. Marocico, A. Emin Tatar

Center for Information Technology
University of Groningen

Friday, July 4th 2025

# Outline

# Introduction to Classification

## Classification

Definition

# Introduction to Classification

## Classification                                                            Definition

Classification is a type of supervised learning where the model learns from labeled data to predict the class of new observations based on past data.

# Introduction to Classification

## Classification
Definition

Classification is a type of supervised learning where the model learns from labeled data to predict the class of new observations based on past data.

Classification vs. Regression is a key distinction in supervised learning:

- In classification, the target variable is categorical (e.g., "spam" or "not spam").

# Introduction to Classification

## Classification
Definition

Classification is a type of supervised learning where the model learns from labeled data to predict the class of new observations based on past data.

Classification vs. Regression is a key distinction in supervised learning:

- In classification, the target variable is categorical (e.g., "spam" or "not spam").
- In regression, the target variable is continuous (e.g., predicting a price).

# Classification Types

## Classification Types
Definition

# Classification Types

## Classification Types

Definition

Classification can be broadly divided into two types:

- Binary Classification: The target variable has two classes (e.g., "yes" or "no", "spam" or "not spam"). Numerically, this can always be represented as 0 and 1.

# Classification Types

## Classification Types
Definition

Classification can be broadly divided into two types:

- Binary Classification: The target variable has two classes (e.g., "yes" or "no", "spam" or "not spam"). Numerically, this can always be represented as 0 and 1.

- Multiclass Classification: The target variable has more than two classes (e.g., "cat", "dog", "weasel"). In this case, the model predicts one of several possible categories.

# Classification Algorithms

## Classification Algorithms                                                    Definition

# Classification Algorithms

## Classification Algorithms
Definition

Classification algorithms are designed to learn from labeled data and make predictions about the class of new, unseen data. Some common algorithms include:

- Logistic Regression: Despite its name, it is used for binary classification. It models the probability that a given input belongs to a particular class.

# Classification Algorithms

## Classification Algorithms
Definition

Classification algorithms are designed to learn from labeled data and make predictions about the class of new, unseen data. Some common algorithms include:

- Logistic Regression: Despite its name, it is used for binary classification. It models the probability that a given input belongs to a particular class.
- k-Nearest Neighbors (k-NN): A non-parametric method that classifies a data point based on the classes of its nearest neighbors in the feature space.

# Classification Algorithms

## Classification Algorithms
Definition

Classification algorithms are designed to learn from labeled data and make predictions about the class of new, unseen data. Some common algorithms include:

- Logistic Regression: Despite its name, it is used for binary classification. It models the probability that a given input belongs to a particular class.
- k-Nearest Neighbors (k-NN): A non-parametric method that classifies a data point based on the classes of its nearest neighbors in the feature space.
- Decision Trees: A tree-like model that splits the data into subsets based on feature values, leading to a decision about the class label.

# Classification Algorithms

## Classification Algorithms
Definition

Classification algorithms are designed to learn from labeled data and make predictions about the class of new, unseen data. Some common algorithms include:

- Logistic Regression: Despite its name, it is used for binary classification. It models the probability that a given input belongs to a particular class.
- k-Nearest Neighbors (k-NN): A non-parametric method that classifies a data point based on the classes of its nearest neighbors in the feature space.
- Decision Trees: A tree-like model that splits the data into subsets based on feature values, leading to a decision about the class label.
- Support Vector Machines (SVM): A method that finds the hyperplane that best separates the classes in the feature space.

# Classification Algorithms

## Classification Algorithms                                                      Definition

Classification algorithms are designed to learn from labeled data and make predictions about the class of new, unseen data. Some common algorithms include:

- Logistic Regression: Despite its name, it is used for binary classification. It models the probability that a given input belongs to a particular class.
- k-Nearest Neighbors (k-NN): A non-parametric method that classifies a data point based on the classes of its nearest neighbors in the feature space.
- Decision Trees: A tree-like model that splits the data into subsets based on feature values, leading to a decision about the class label.
- Support Vector Machines (SVM): A method that finds the hyperplane that best separates the classes in the feature space.
- More advanced algorithms like Random Forests, Gradient Boosting, and Neural Networks.

# Logistic Regression

## Logistic Regression

Definition

# Logistic Regression

## Logistic Regression
Definition

Logistic Regression models the probability that the target variable $y$ belongs to a particular class. The logistic function (sigmoid) is used to map predicted values to probabilities between 0 and 1. The decision boundary is determined by the threshold (commonly 0.5) for classifying observations into different classes.

# Logistic Regression

## Logistic Regression

Logistic Regression models the probability that the target variable $y$ belongs to a particular class. The logistic function (sigmoid) is used to map predicted values to probabilities between 0 and 1. The decision boundary is determined by the threshold (commonly 0.5) for classifying observations into different classes.

The sigmoid function is defined as:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

where $z$ is a linear combination of the input features.