# Machine learning

Courses: Industrial automation, Automatizace pro průmyslovou praxi

CTU, FS, U12110

Matouš Cejnek

# Contents

- Introduction
- Classification & Clustering
- Cross-validation

# Introduction

# Machine learning (ML)

A branch of artificial intelligence (AI) where **computers learn from data to improve their performance** on a task.

Or in other words:

The **ability of computer systems to learn patterns** and make decisions without explicit programming.

# Machine learning utilization

- **Data analysis**:
  - Pre-processing
  - Core-analysis
  - Post-processing
- Applied data science
  - Data product deployment in production
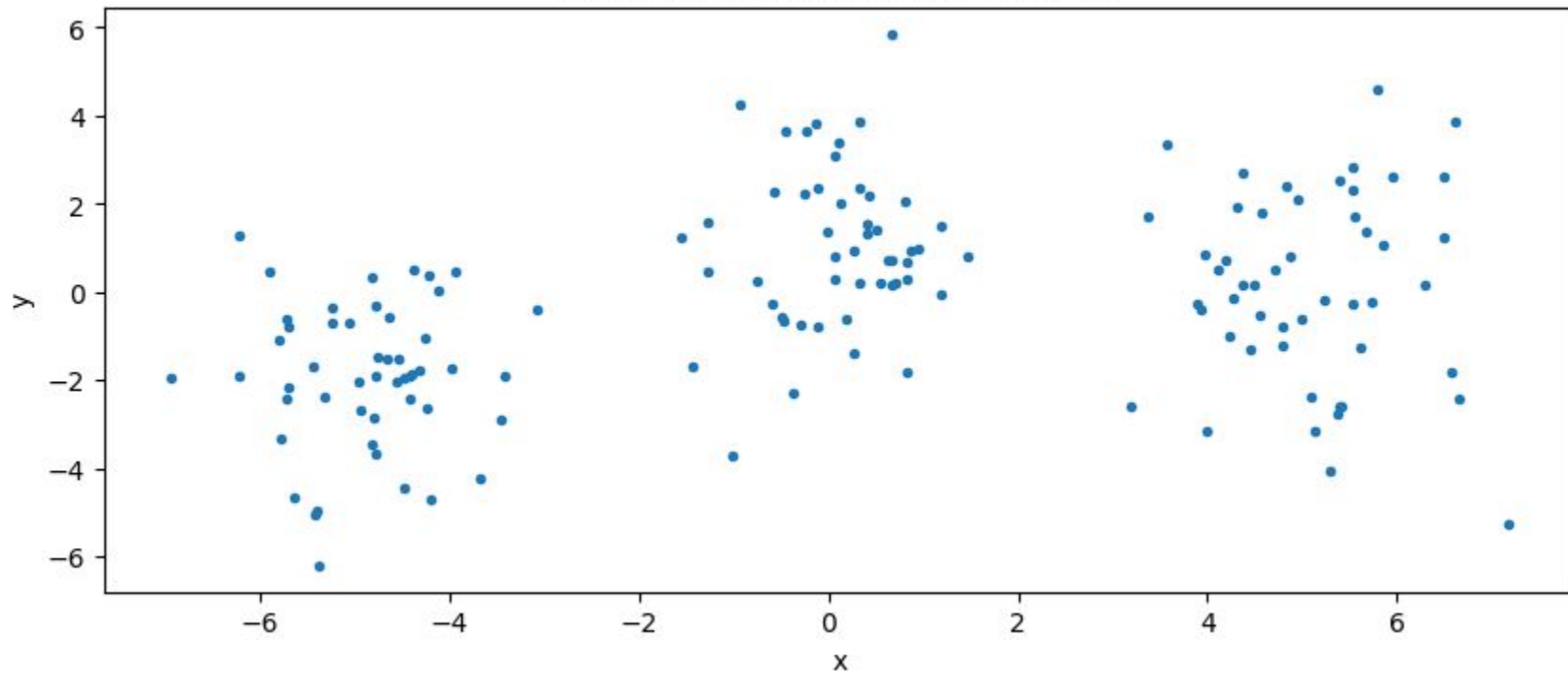
# ML data analysis

Most common:

- Regression
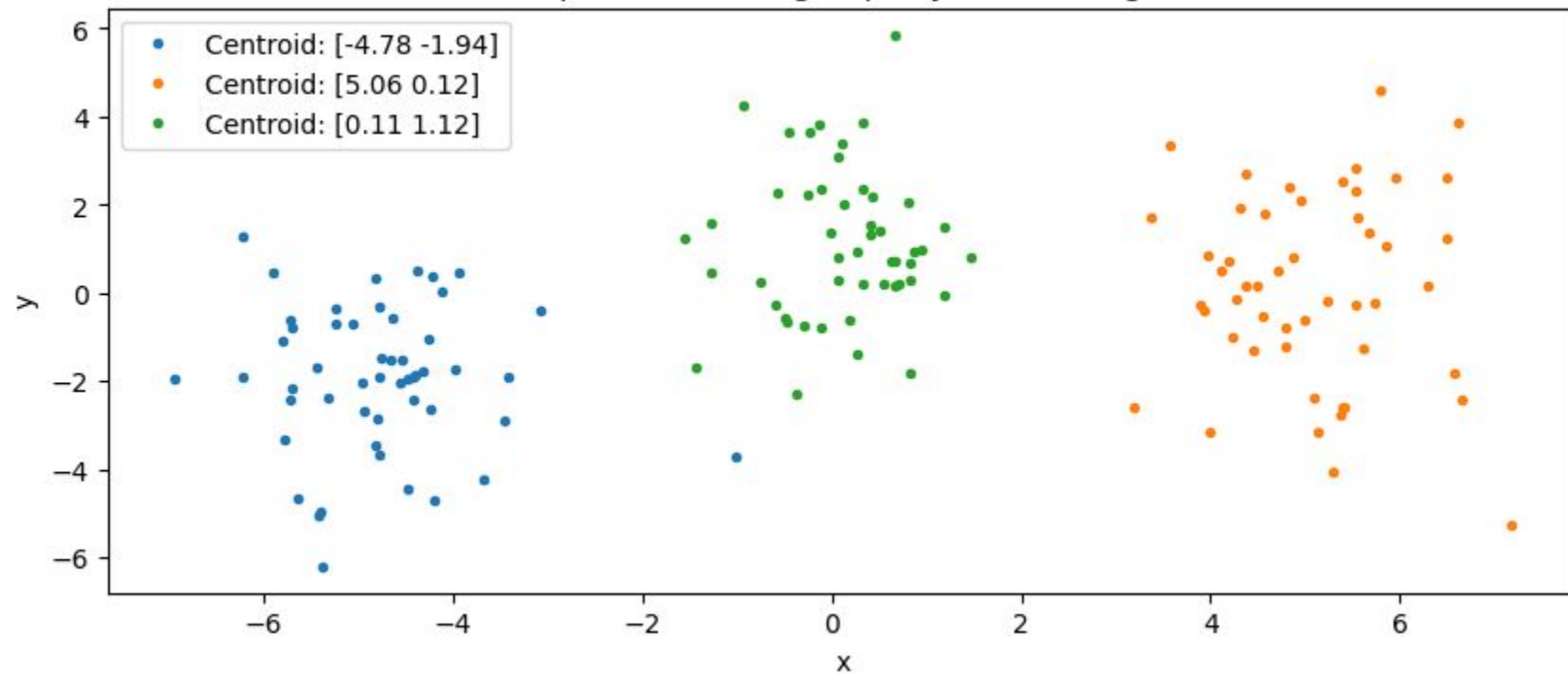- **Classification**
- **Clustering**

# Clustering

Clustering is an **unsupervised** learning technique in machine learning that involves grouping similar data points into distinct subsets or clusters.

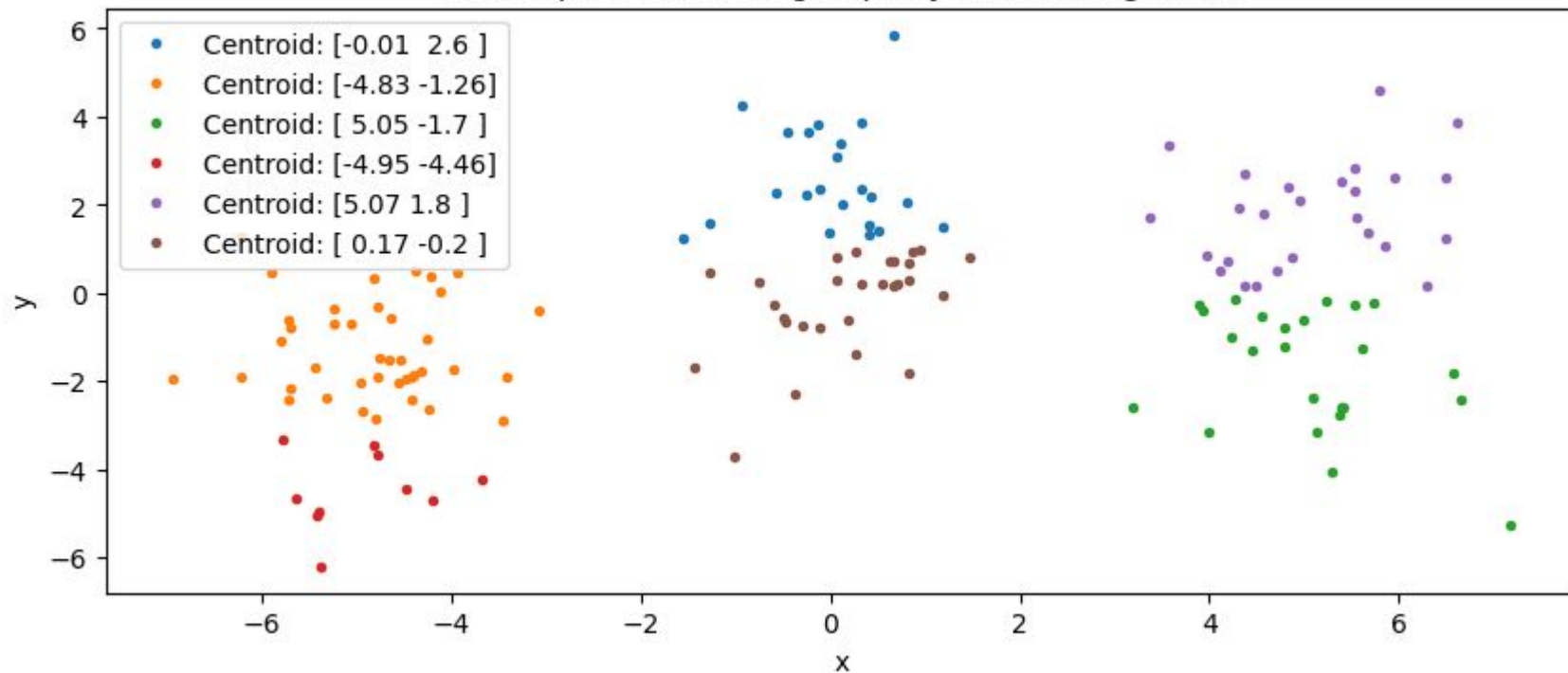- We're unsure about cluster shapes or their quantity

Unknown data with visible clusters

Data separated into 3 groups by Kmeans algorithm
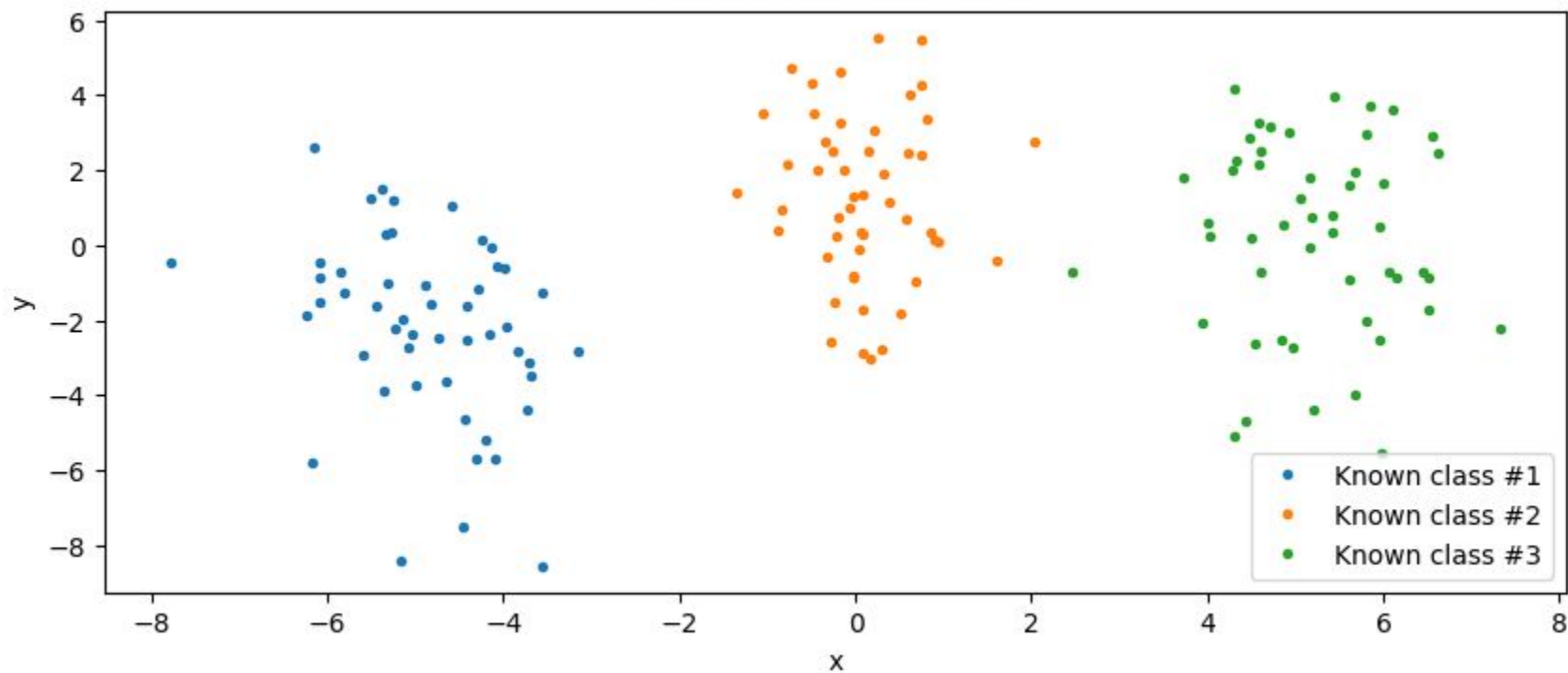
# Popular clustering algorithms

- Kmeans
- DBscan
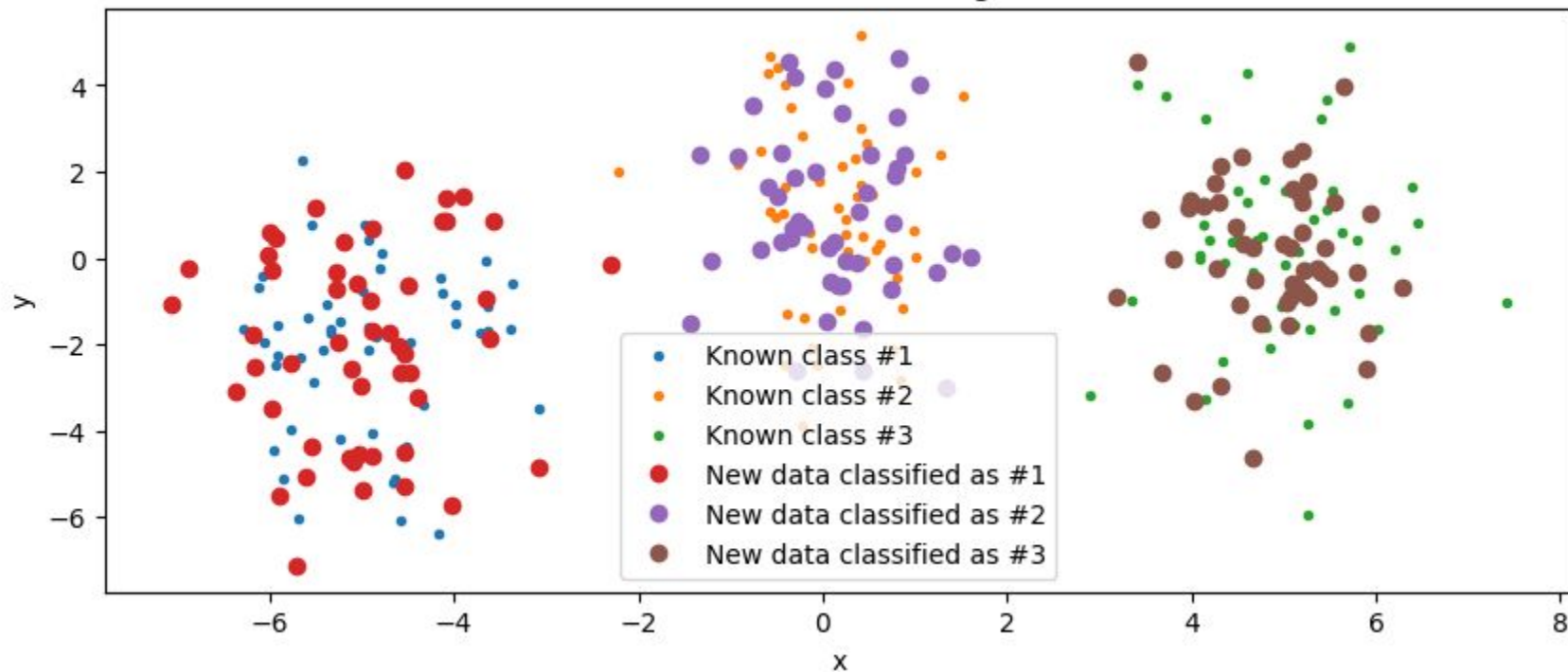- Self organizing map (SOM)
- …

# Classification

Classification is a **supervised** learning task in machine learning where the goal is to assign predefined labels or categories to input data based on its features.

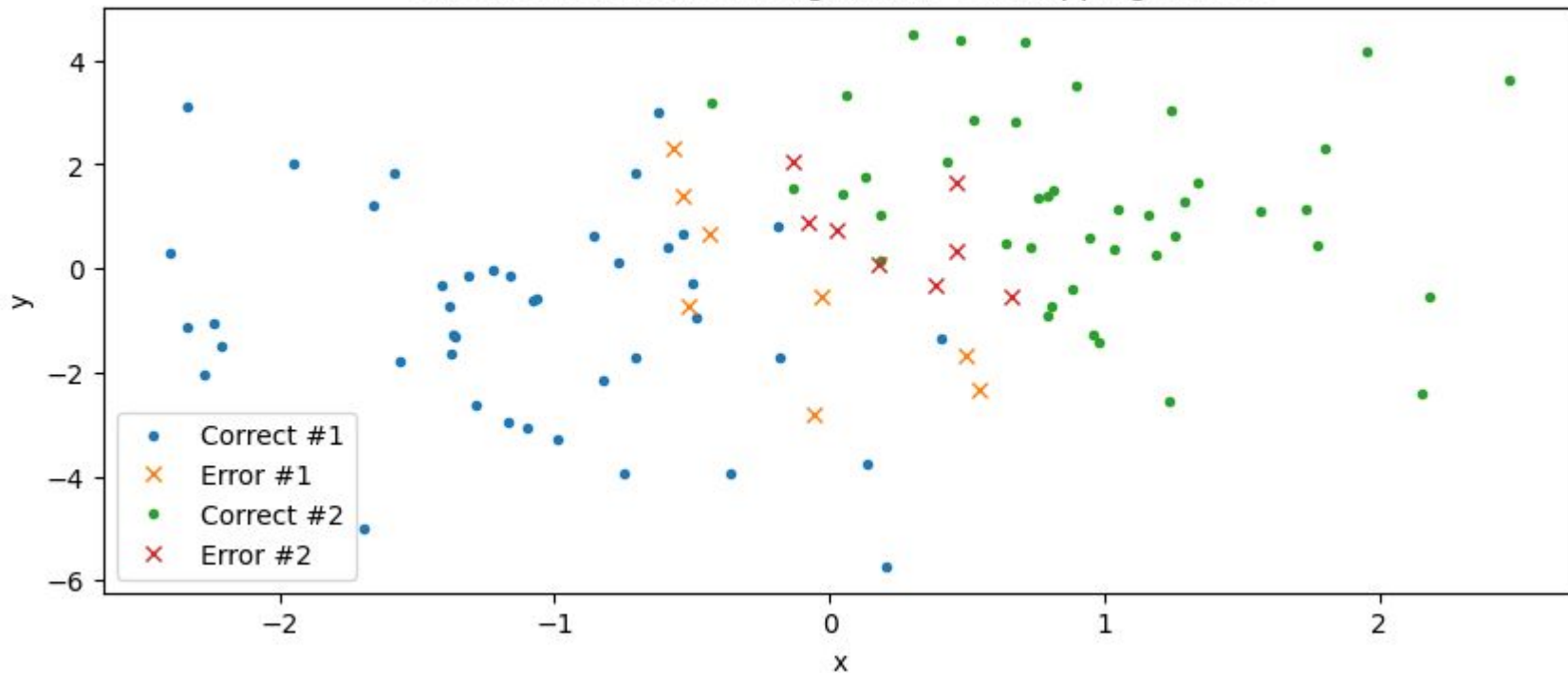- We have defined classes and clear expectations for each.

Classification with SVM algorithm.

Classification with SVM algorithm - overlapping classes

# Popular classification algorithms

- K-Nearest Neighbors
- Support vector machines (SVM)
- Decision trees
- …

# Clustering vs classification

- **Classification:** Assigns predefined labels to data points based on their inherent properties.
- **Clustering:** Seeks optimal label groupings within data, uncovering natural patterns.

# Clustering vs classification

- **Classification:** We have defined classes and clear expectations for each.
- **Clustering:** We're unsure about cluster shapes or their quantity

# Classifier evaluation

# Confusion matrix

Confusion matrix is a specific table layout that allows visualization of the performance of an algorithm.

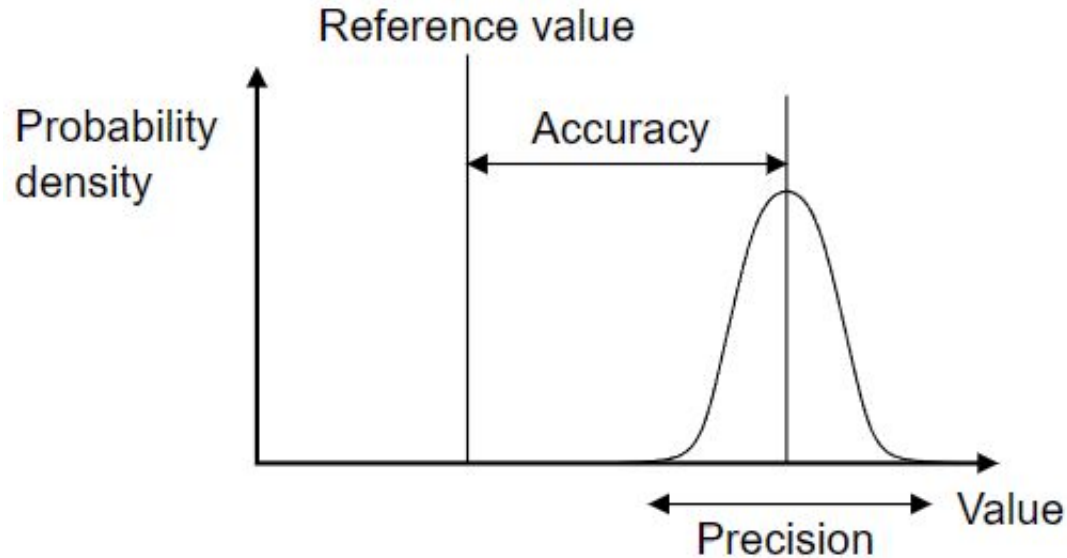| | | Predicted condition | |
|---|---|---|---|
| | Total population (P+ N) | Positive (PP) | Negative (PN) |
| Actual condition | Positive (P) | True positive (TP) | False negative (FN) |
| | Negative (N) | False positive (FP) | True negative (TN) |

# Accuracy and Precision

- **Accuracy:** $$ACC = \frac{TP + TN}{P + N}$$

- **Precision:** $$PPV = \frac{TP}{TP + FP}$$

# Precision vs Accuracy

# Sensitivity and specificity

- **Sensitivity:**

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$
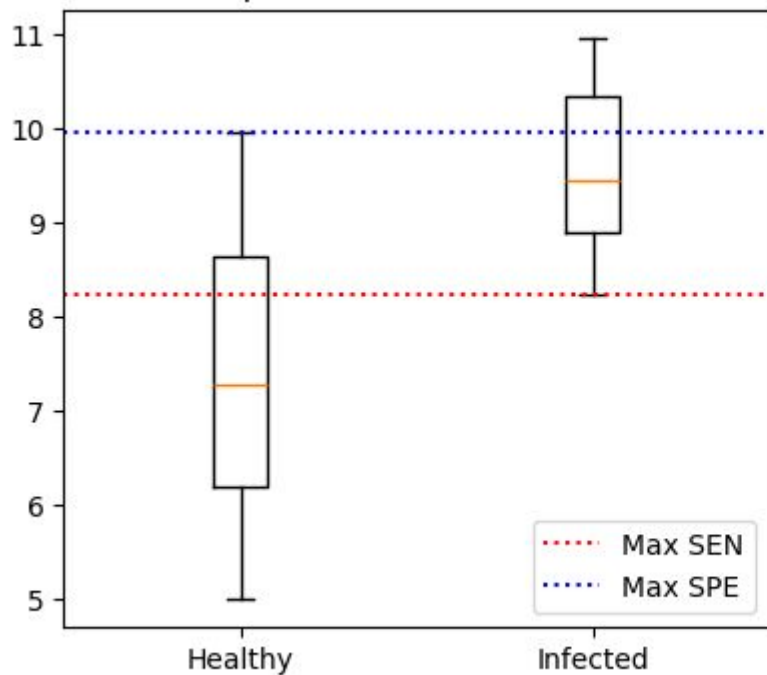
  Also known as: Recall, Hit rate, True positive rate

- **Specificity:**
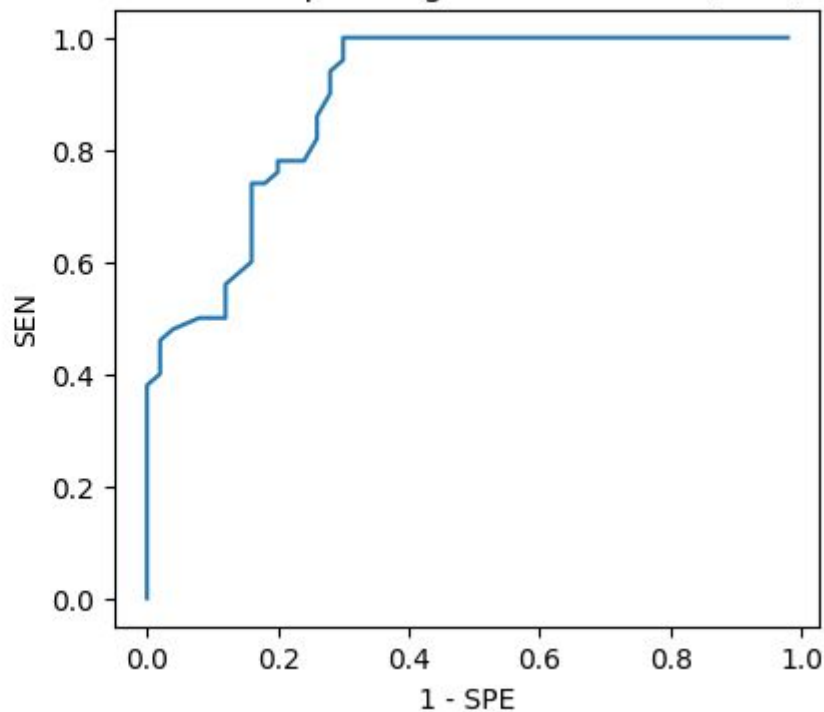
$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

  Also known as: Selectivity, True negative rate

# Receiver Operating Characteristics (ROC)

- ROC curve, is a graphical plot that illustrates the performance of a classifier at varying threshold values.
- **AUROC** (area under the ROC curve) is very useful metric

# Cross-validation

# Validation

- Validation is a method for evaluating a predictive model performance by splitting the dataset into training and testing sets.
- In the ideal scenario, the prediction error on the training set should ideally be equal to or smaller than the training set error.

# Cross-validation

Cross-validation is a method for evaluating a predictive model performance by systematically splitting the dataset into training and testing sets multiple times.

- N-fold cross-validation
- Leave-p-out cross-validation

# Cross-validation

Cross-validation is a method for evaluating a predictive model performance by systematically splitting the dataset into training and testing sets multiple times.

- N-fold cross-validation
- Leave-p-out cross-validation

# Cross-validation

**N-fold cross-validation**:

You train your model **n** times, each time using a different fold as the test set and the remaining folds as the training set.

**Leave-p-out cross-validation**:

You train your model multiple times, each time leaving out just **p**-number of data point for testing.

# N-fold cross-validation

1. **Split data** into n-folds
2. Use **one fold for testing** and other folds for model training
3. **Repeat step 2** with different folds (train from scratch)
4. **Average the results** over all tested folds
- Folds are chosen at the beginning.
- Folds are not overlapping.

# Leave-p-out cross-validation

1. **Select p** number of **samples** from data
2. **Use selected samples for testing**, other data for training
3. **Repeat steps 1 and 2** (train from scratch)
4. **Average the results** over all tests

# Overtraining

# Overtraining (overfitting)

Overtraining, occurs when a model learns the training data too well, including its noise and outliers, to the extent that it negatively impacts the model's performance on new, unseen data.

- Training dataset do not represent the process fully
- We need models that can generalize