# Lecture 13

⟹ Intro to Scikit — Learn

MACHINE LEARNING

labels

SUPERVISED LEARNING

no labels

UNSUPERVISED LEARNING

classification (discrete)

regression (continuous)

clustering

density estimation

dimensionality reduction

Terminology

* Data recorded in matrix form

$$\hookrightarrow [X] = [N_{samples}, N_{features}]$$

# datapoints    # attributes

* Some data will have labels — stars, galaxies etc.
(sklearn calls these "TARGETS")

- Scikit-Learn Workflow
  1. Instantiate an estimator object
  2. Fit estimator on data and labels
  3. Predict new labels
  4. Find model parameters

  ⟹ Usually partition dataset into TRAINING and TESTING sets.

---

- Supervised Learning — labels

  * Classification algorithms are trained on labelled data, and used to classify new object features.

  * Regression (or "fitting") is the continuous form of classification.

---

- Unsupervised Learning — no labels.

  * Use the data to discover its own labels.

  * Clustering (group similar data), density estimation (find the PDF), dimensionality reduction (find important features).