# Day 02 - Classification

# SDSS Data Set

# SDSS Data Set

# SDSS Data Set



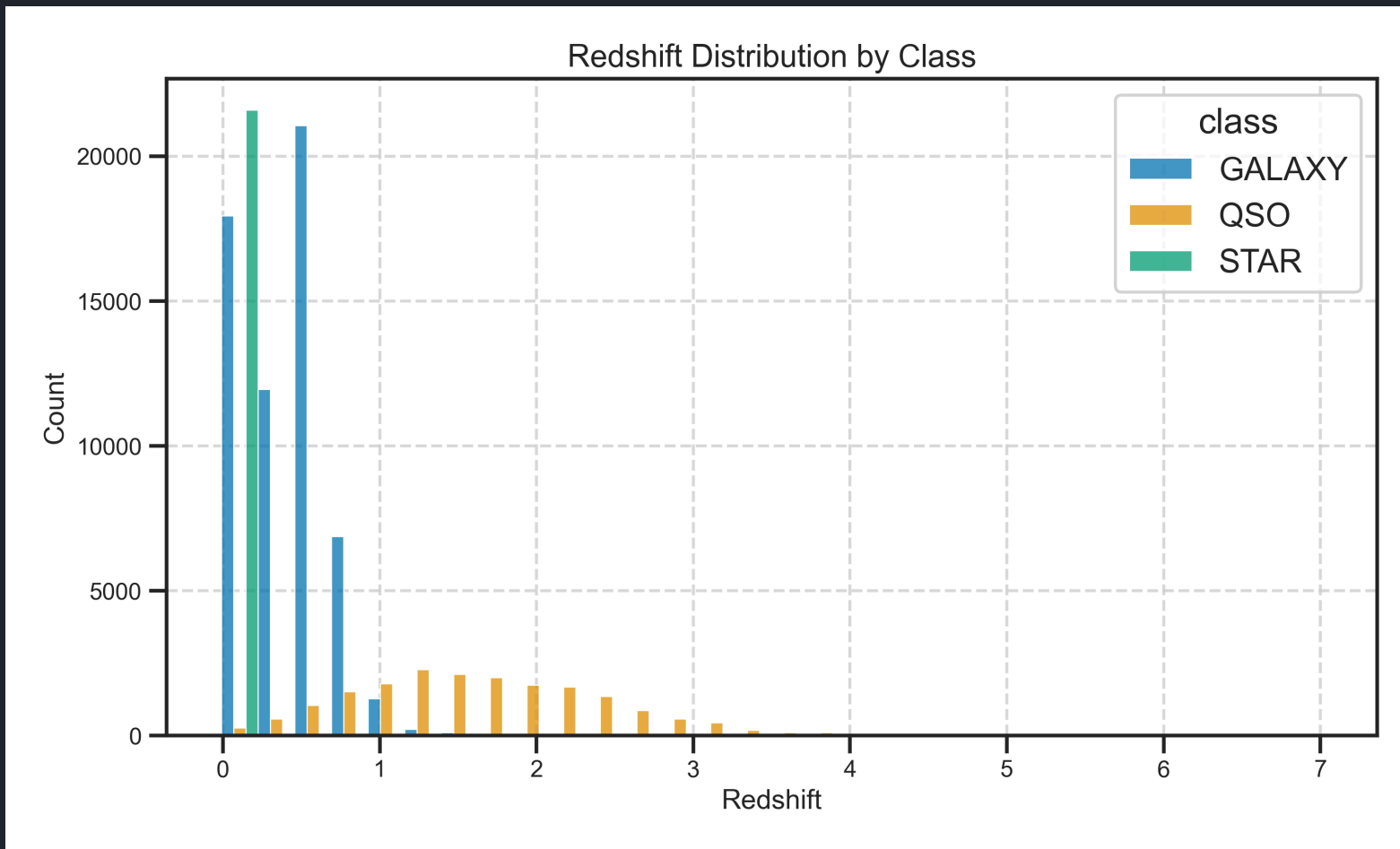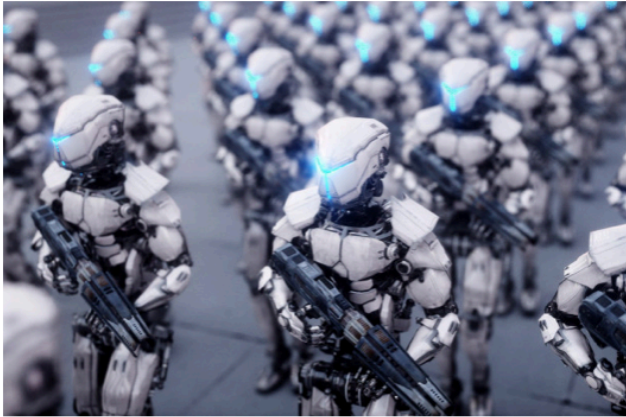Redshift Distribution by Class

# Classification Task

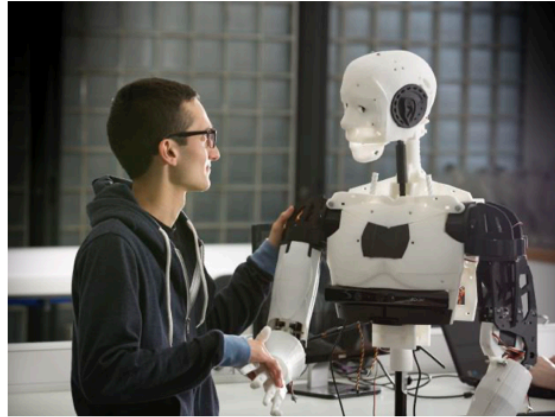- Using the SDSS data set, we will classify objects as stars or quasars.
  - At first, we will only use the color information (u-g, g-r, r-i) to classify objects.
  - Later, we will add the redshift information (z) to improve our classification.
- Then, we will perform a 3-class classification to distinguish between stars, quasars, and galaxies; here we will use all available features including redshift.

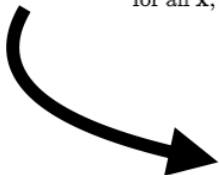# Machine Learning



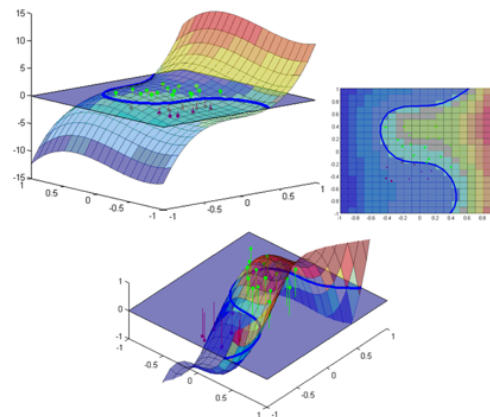What society thinks we do



What our friends think we do



What our families thinks we do



What my boss thinks we do



What we think we do



What we actually do

# Sci-Kit-Learn Classification

- Sci-Kit-Learn is a powerful Python library for machine learning.

- It provides a wide range of classification algorithms, including:
    - k-Nearest Neighbors (kNN) & Logistic Regression
    - Decision Trees & Random Forests
    - Support Vector Machines (SVM)

- It also includes tools for model evaluation, such as cross-validation and confusion matrices.

https://scikit-learn.org/stable/index.html

Simplified Approach

# K-Nearest Neighbors (kNN)

- kNN is a simple and intuitive classification algorithm.

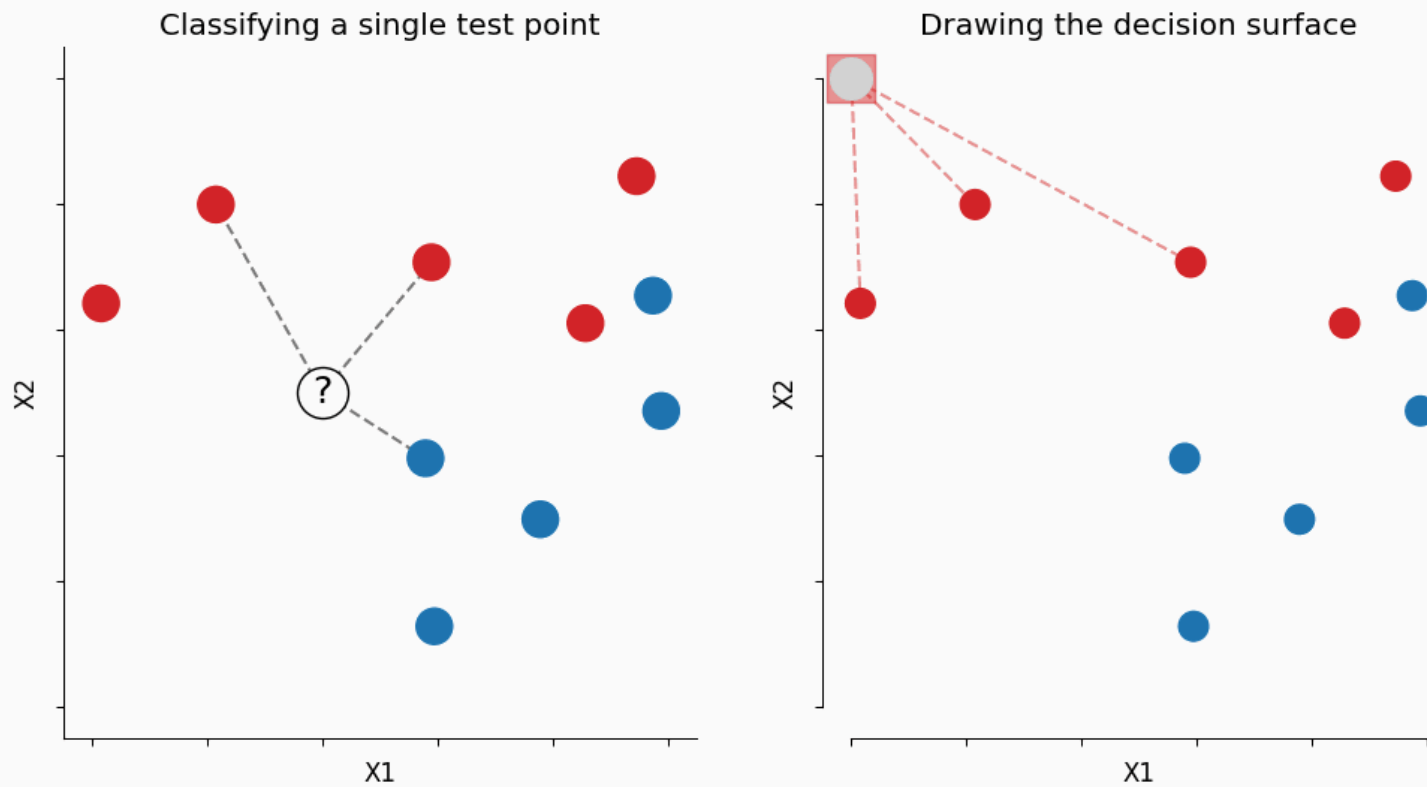- It classifies a data point based on the majority class of its k nearest neighbors in the feature space.

- The distance metric (e.g., Euclidean distance) is used to determine the nearest neighbors.

- kNN is a non-parametric method, meaning it makes no assumptions about the underlying data distribution.

- It is sensitive to the choice of k and the distance metric.

# K-Nearest Neighbors (kNN)

# Today's Activity

- We will implement a kNN classifier using Sci-Kit-Learn to classify stars and quasars from the SDSS data set.
- We will:
    - i. Load the SDSS data set and preprocess it.
    - ii. Split the data into training and testing sets.
    - iii. Train a kNN classifier on the training set.
    - iv. Evaluate the classifier's performance on the testing set.
- We will focus on the evaluation metrics and visualizations to understand the classifier's performance.