

# Different Types Of Datasets Splitting Methods

There are several methods for splitting datasets for machine learning and data analysis. Each method serves different purposes and has its advantages and disadvantages. Here are some common datasets splitting methods, along with step-by-step explanations for each:

## 1. **Train-Test Split (Holdout Method):**

- Purpose: To create two separate sets, one for training and one for testing the model.
- Steps:
  1. Randomly shuffle the dataset to ensure the data is well-distributed.
  2. Split the data into two parts, typically with a ratio like 70-30 or 80-20, where one part is for training and the other for testing.
  3. Train your machine learning model on the training set.
  4. Evaluate the model's performance on the test set.

## 2. **K-Fold Cross-Validation:**

- Purpose: To assess the model's performance by training and testing it on different subsets of the data.
- Steps:
  1. Divide the dataset into k equal-sized folds.
  2. For each fold (1 to k), treat it as a test set, and the remaining (k-1) folds as the training set.
  3. Train and evaluate the model on each of the k iterations.
  4. Calculate performance metrics (e.g., accuracy) by averaging the results from all iterations.

## 3. **Stratified Sampling:**

- Purpose: To ensure that the proportion of different classes in the dataset is maintained in the train and test sets.
- Steps:
  1. Identify the target variable (e.g., class labels).
  2. Stratify the data by the target variable to create representative subsets.
  3. Perform a train-test split on these stratified subsets to maintain class balance in both sets.

## 4. **Time Series Split:**

- Purpose: For time series data, where the order of data points matters.
- Steps:
  1. Sort the dataset based on the time or date variable.
  2. Divide the data into training and testing sets such that the training set consists of past data, and the testing set contains future data.

## 5. **Leave-One-Out Cross-Validation (LOOCV):**

- Purpose: To leave out a single data point as the test set in each iteration.
- Steps:
  1. For each data point in the dataset, create a training set with all other data points.
  2. Train and test the model for each data point separately.
  3. Calculate the performance metrics based on the predictions from each iteration.

## 6. **Group K-Fold Cross-Validation:**

<https://aiquest.org/>

<https://github.com/ahammadmejbah>

## Different Types Of Datasets Splitting Methods

- Purpose: To account for groups or clusters in the data.
- Steps:
  1. Identify the group or cluster variable (e.g., patient ID).
  2. Divide the data into folds while ensuring that each fold contains all data points from a specific group.
  3. Train and test the model using group-based cross-validation.

### 7. **Bootstrapping:**

- Purpose: To create multiple subsets of the data with replacement for estimating model performance.
- Steps:
  1. Randomly sample data points with replacement to create multiple bootstrap samples.
  2. Train and evaluate the model on each bootstrap sample.
  3. Calculate performance metrics based on the results of each sample.

The choice of dataset splitting method depends on the specific problem, data characteristics, and the goal of your analysis. It's important to select an appropriate method to ensure reliable model evaluation and generalization.