



**k-Fold Cross-Validation** → Is a technique that minimizes the disadvantages of hold-out method. k-Fold introduces a new way of splitting the dataset which helps to overcome the "test only once bottleneck".

k-Fold gives a more stable and trustworthy result as previous method since training and testing is performed on several different parts of the dataset.

Still, k-Fold method has a disadvantage. Increasing  $k$  results in training more models and the training process may be really expensive and time-consuming.

### Algorithm:

1. Pick a number of folds -  $k$ . Usually,  $k$  is 5 to 10 but you can choose any number which is less than the dataset's length.
2. Split the dataset into  $k$  equal (if possible) parts (they are called folds)
3. Choose  $k-1$  folds which will be the training set. The remaining fold will be the test set.
4. Train the model on the training set. On each iteration of cross-validation, you must train a new model independently of the model trained in the previous iteration.
5. Validate on the test set.
6. Save the result of the validation.
7. Repeat steps 3-6  $k$  times. Each time use the remaining fold as the test set. In the end, you should have validated the model on every fold that you have.
8. To get the final score average the results that you got on step 6.

### Implementation:

```
import numpy as np
from sklearn.model_selection import KFold
```

```
X = np.array([[1,2], [3,4], [1,2], [3,4]])
```

```
y = np.array([1,2,3,4])
```

```
kf = KFold(n_splits=2)
```

```
for train_index, test_index in kf.split(x):
    print("TRAIN: ", train_index, "TEST: ", test_index)
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]
```

**Leave-one-out cross-validation (LOOCV)** → Is an extreme case of k-Fold CV where  $k$  is equal to  $n$ , where  $n$  is the number of samples in the dataset.

The greatest advantage of LOOCV is that it doesn't waste much data. LOOCV is more computationally expensive than k-Fold, it make take plenty of time to cross-validate the model using LOOCV. Thus, the Data Science community has a general rule based on empirical evidence and different researches, which suggest that 5 or 10 fold cross-validation should be preferred over LOOCV.

### Implementation:

```
import numpy as np
from sklearn.model_selection import LeaveOneOut
```

```
X = np.array([[1, 2], [3, 4]])
```

```
y = np.array([1, 2])
```

```
loo = LeaveOneOut(1)
```

```
for train_index, test_index in loo.split(X):
```

```
    print("TRAIN:", train_index, "TEST:", test_index)
```

```
    X_train, X_test = X[train_index], X[test_index]
```

```
    y_train, y_test = y[train_index], y[test_index]
```

**Leave-p-out cross-validation (LpOCV)** → Is similar to LOOCV as it creates all the possible training and test sets by using  $p$  samples as the test set.

LpOCV has all disadvantages of the LOOCV, but, nevertheless, it's as robust as LOOCV

### Algorithm:

1. Choose  $p$  samples from the dataset which will be the test set.
2. The remaining  $n-p$  samples will be the training set.
3. Train the model on the training set. On each iteration, a new model must be trained.
4. Validate on the test set.
5. Save the result of the validation.
6. Repeat steps 2-5  $\binom{n}{p}$  times.
7. To get the final score average the results that you got on step 5.