





/ Validation Strategies

SimpleSplit, CrossValidation, Stratification,...



/ Validation









Is the process of testing that the model performs well on new unseen data. We should always validate our models.









Supervised and Unsupervised Data

Supervised data

Labelled data		Labelled data	
	Dog		18 lbs
	Dog		14 lbs
	Cat		12 lbs
	Cat		9 lbs

Unsupervised data

Unlabelled data	
	
	



Simple Split Validation

- Dataset (usually supervised)
 - **Train** The data that the model learns
 - **Validation** Data to check that the model performs well
- New data (usually unsupervised)
 - **Test** The new data to do inference

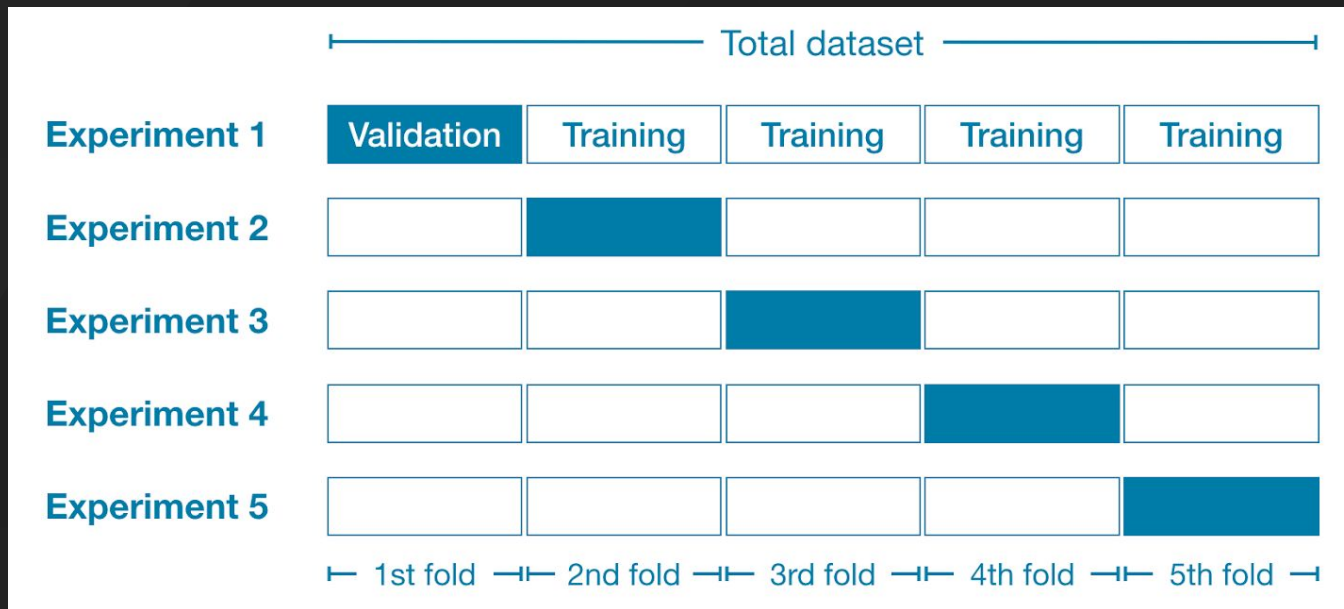




Cross Validation (aka K-Fold)

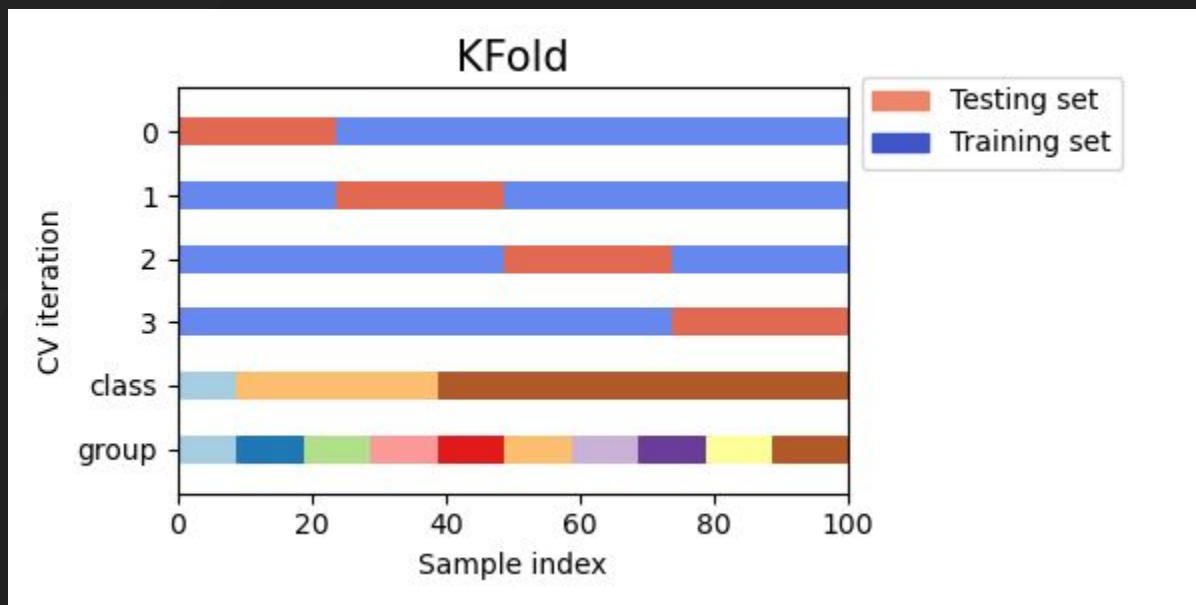
/ Cross Validation with 5 Folds (K=5). This validates more! (5 times)

/ Remember: On each experiment a new model is trained from scratch.



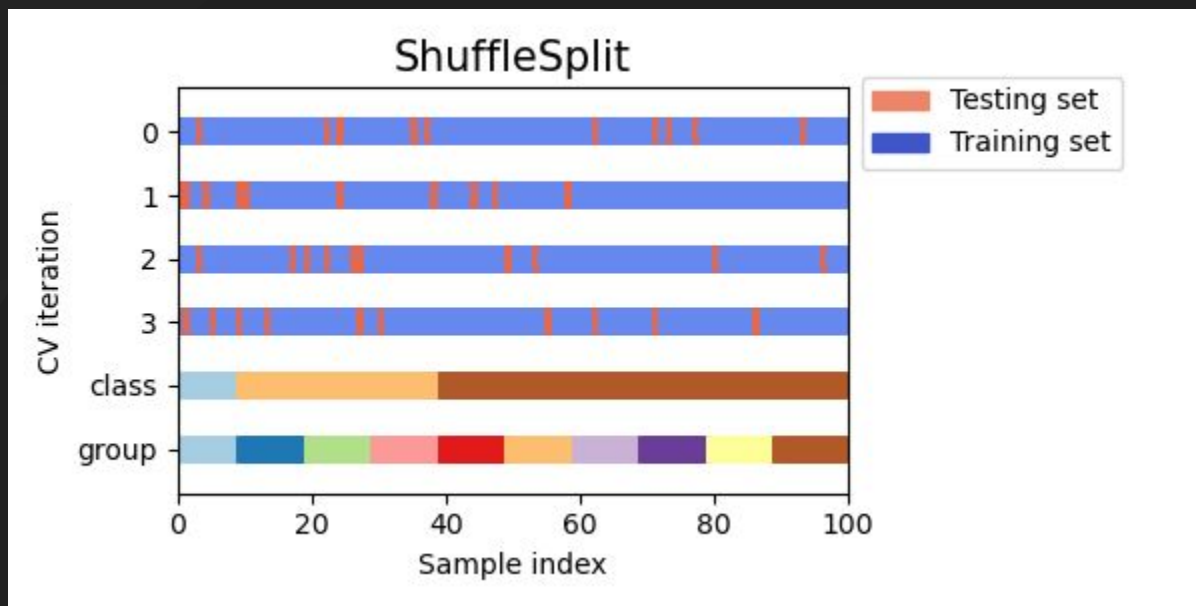


Cross Validation





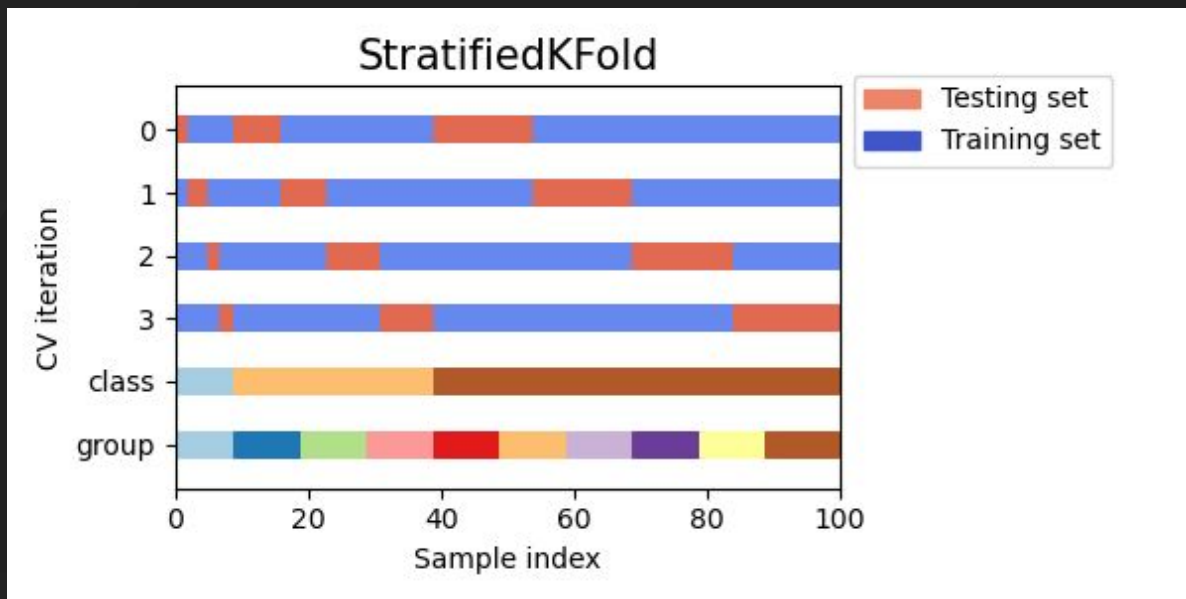
Cross Validation randomized





Cross Validation with Stratification

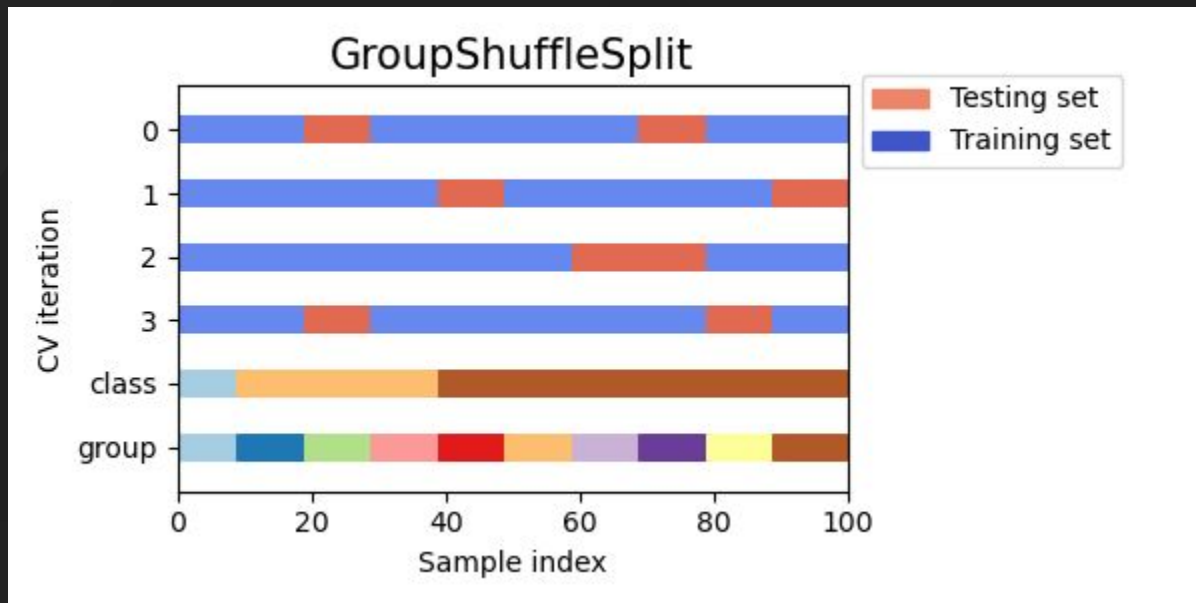
/ Each fold have the same class distribution. **Very useful for classification!**





Cross Validation with Groups

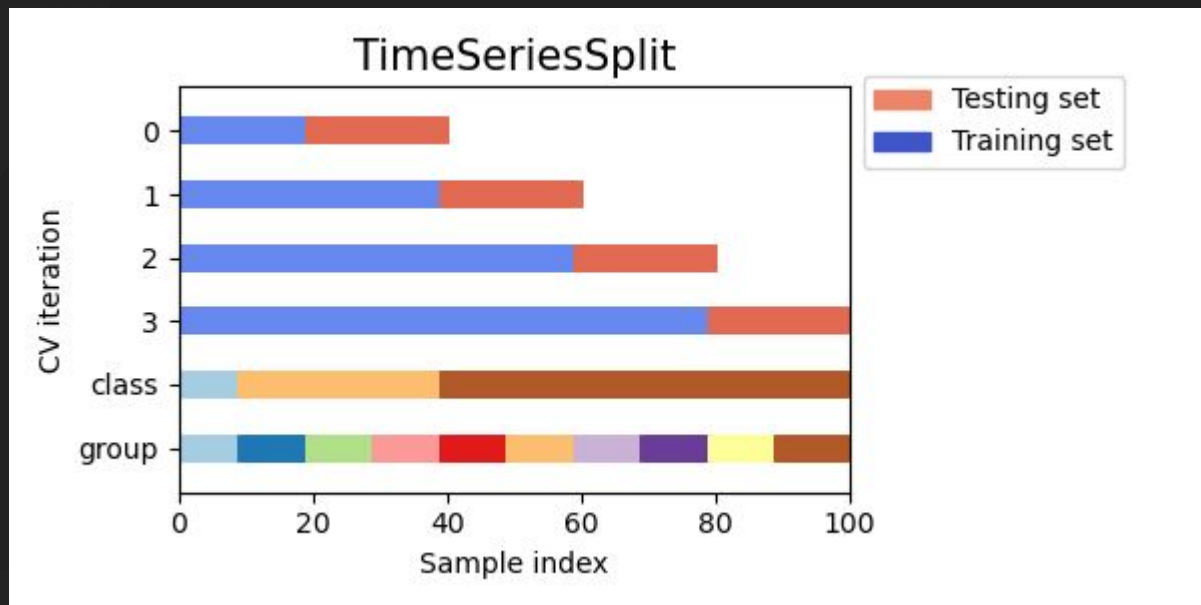
/ Useful when same samples are from the same individual.
E.g. patient have 3 records, and we want to validate on new patients.






Cross Validation for Time Series

/ Validate always on the future data





Sklearn validation

/ Many options to choose

`model_selection.GroupKFold([n_splits])`

`model_selection.GroupShuffleSplit(...)`

`model_selection.KFold([n_splits, shuffle, ...])`

`model_selection.LeaveOneGroupOut()`

`model_selection.LeavePGroupsOut(n_groups)`

`model_selection.LeaveOneOut()`

`model_selection.LeavePOut(p)`

`model_selection.PredefinedSplit(test_fold)`

`model_selection.RepeatedKFold(*[, n_splits, ...])`

`model_selection.RepeatedStratifiedKFold(*[, ...])`

`model_selection.ShuffleSplit([n_splits, ...])`

`model_selection.StratifiedKFold([n_splits, ...])`

`model_selection.StratifiedShuffleSplit(...)`

`model_selection.TimeSeriesSplit([n_splits, ...])`

K-fold iterator variant with non-overlapping groups.

Shuffle-Group(s)-Out cross-validation iterator

K-Folds cross-validator

Leave One Group Out cross-validator

Leave P Group(s) Out cross-validator

Leave-One-Out cross-validator

Leave-P-Out cross-validator

Predefined split cross-validator

Repeated K-Fold cross validator.

Repeated Stratified K-Fold cross validator.

Random permutation cross-validator

Stratified K-Folds cross-validator.

Stratified ShuffleSplit cross-validator




Time Series cross-validator



Validation for multilabel classification

/ [iterative-stratification](#) is a project that provides scikit-learn compatible cross validators with stratification for multilabel data.

- MultilabelStratifiedKFold
- MultilabelRepeatedStratifiedKFold
- MultilabelStratifiedShuffleSplit

Samples		
		
Labels (t)		
[1 0 1]	[0 1 0]	[1 1 1]

MultiLabel dataset



/ Q&A

What are your doubts?

