

Common Pitfalls and recommended practices

Sneaking, scikit-learn data preprocessing and uncertainty are covered in this chapter. Inconsistent preprocessing causes model failure since train and test data should be scaled identically. Pipeline automation and dataloss minimization may overstate model training results. Choose ~~features~~ features from training data, not test. Managing uncertainty is hard. An integer for 'random-state' produces repeatable results, although cross-validation may be unreliable. Randomstate instances test models in random situations to make them resilient. RandomForest and cross validation splitter yield mixed fold results based on "random-state". Avoid 'random-state=None' and follow cv unpredictability and RandomForestClassifier or other random estimator distribution guidelines.