# LESSON 3: End-to-end ML

## CARSTEN EIE FRIGAARD

AUTUMN 2020



$$\text{Performance} = \frac{1}{5} \sum_{i=1}^{5} \text{Performance}_i$$

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E." — Mitchell (1997).

# Agenda
End-to-end Machine Learning

1. Admin
   - ▶ Zoom undervisnings, lektions videoer på BB,
   - ▶ Afleveringer, grupper, etc.
   - ▶ GITMAL og overskrivning af jeres filer!
2. Indledende undersøgelser og valg af data til slut-projekt (O4).
3. Generel repetition af § 2.
4. Algo. og Model selection, K-fold Cross validation.

# Opg. L03 Beskrivelse af eget slutprojekt.pdf

Dit datasæt fra f.eks. `https://www.kaggle.com...`

(brug min login: user=cef@ase.au.dk, password=test123)

# Opg. L03 Beskrivelse af eget slutprojekt.pdf

…eller UCI https://archive.ics.uci.edu/ml/index.php…

# End-to-end ML

'The Map'

# ML Algorithm Selection and Model Selection

Manually Choosing an Algorithm and Tuning a Model..

- ▶ algorithm selection
  (choose a $h()$).
- ▶ model selection
  (set hyperparameters on $h()$),
- ▶ model evaluation,
- ▶ re-iteration and re-selection!

Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning

Sebastian Raschka
University of Wisconsin-Madison
Department of Statistics
November 2018
sraschka@wisc.edu

**Abstract**

The correct use of model evaluation, model selection, and algorithm selection techniques is vital in academic machine learning research as well as in many industrial settings. This article reviews different techniques that can be used for each of these three subtasks and discusses the main advantages and disadvantages of each techniqu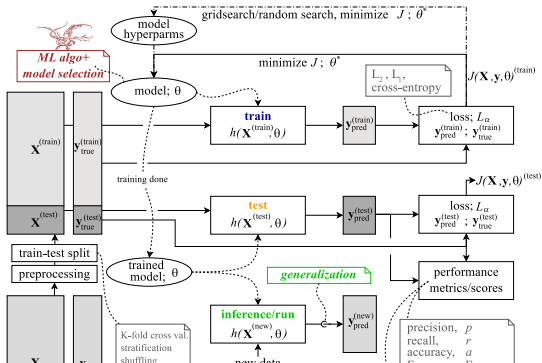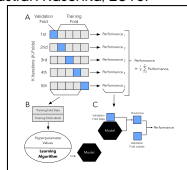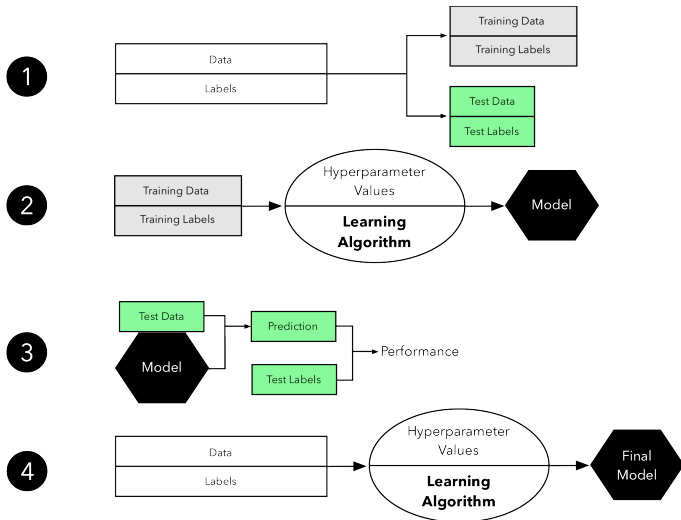e with references to theoretical and empirical studies. Further, recommendations are given to encourage best yet feasible practices in research and applications of machine learning. Common methods such as the holdout method for model evaluation and selection are covered, which are not recommended when working with small datasets. Different flavors of the bootstrap technique are introduced for estimating the uncertainty of performance estimates, as an alternative to confidence intervals via normal approximation if bootstrapping is computationally feasible. Common cross-validation techniques such as leave-one-out cross-validation and k-fold cross-validation are reviewed, the bias-variance trade-off for choosing k is discussed, and practical tips for the optimal choice of k are given based on empirical evidence. Different statistical tests for algorithm comparisons and strategies for dealing with multiple comparisons such as omnibus tests and multiple-comparison corrections are discussed. Finally, alternative methods for algorithm selection, such as the combined F-test 5x2 cross-validation and nested cross-validation, are recommended for comparing machine learning algorithms when datasets are small.

*"Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning",* Sebastian Raschka, 2018.
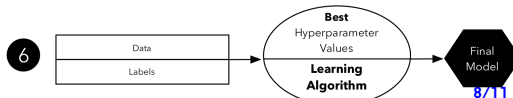
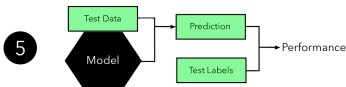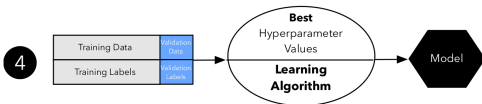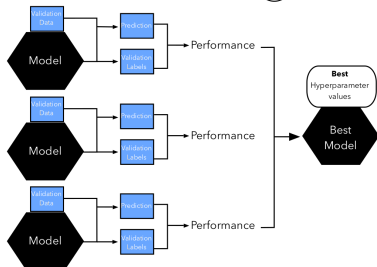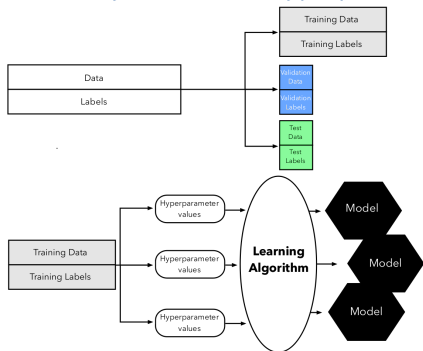# Model Evaluation

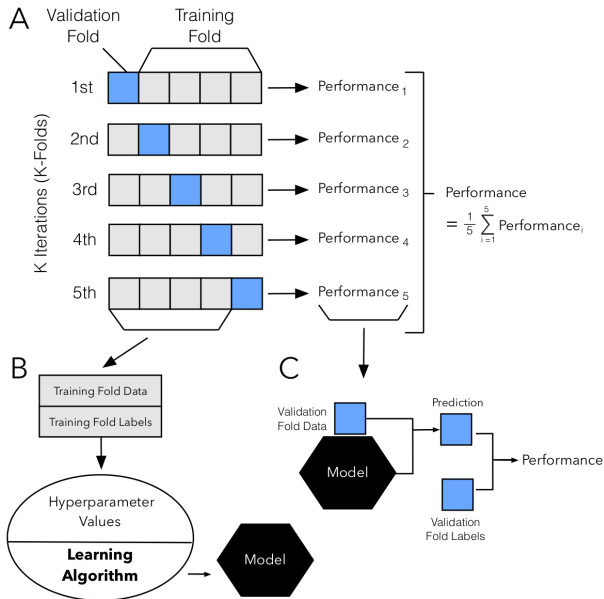## Simple Holdout Method (Train-Test Split)..

# Model Evaluation and Selection

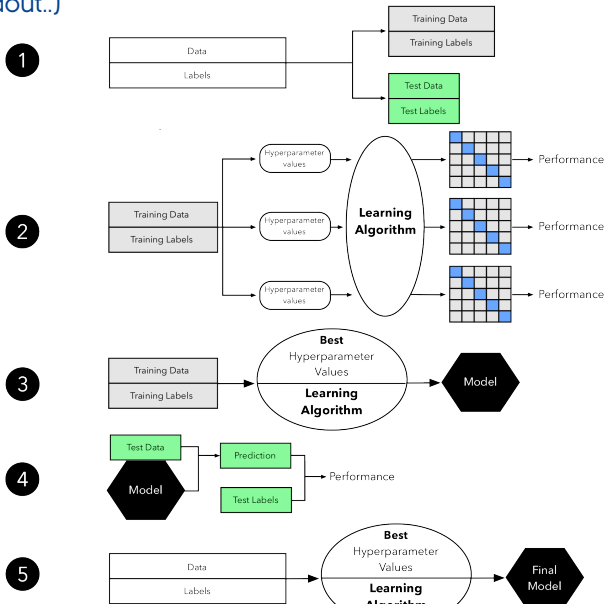Three-way Holdout for Hyperparameter Tuning (Train-Validate-Test Split)...

# Model Evaluation

*k*-fold Cross-Validation Procedure, for *k* = 5..

# Model Evaluation and Selection

*k*-fold Cross-Validation for Hyperparameter Tuning (Somewhat Similar to Treeway Holdout..)

# Scikit-learn K-fold Demo..

Install    User Guide    API    Examples    More ▾

Please **cite us** if you use the software.

sklearn.model_selection.K Fold

Examples using

sklearn.model_selection.KF

## sklearn.model_selection.KFold

*class* sklearn.model_selection.**KFold**(*n_splits=5*, *, *shuffle=False*, *random_state=None*)    [source]

K-Folds cross-validator

Provides train/test indices to split data in train/test sets. Split dataset into k consecutive folds (without shuffling by default).

Each fold is then used once as a validation while the k - 1 remaining folds form the training set.

Read more in the User Guide.

| Parameters: | **n_splits** : *int, default=5* |
| | Number of folds. Must be at least 2. |
| | *Changed in version 0.22:* `n_splits` default value changed from 3 to 5. |
| | **shuffle** : *bool, default=False* |
| | Whether to shuffle the data before splitting into batches. Note that the samples within each split will not be shuffled. |
| | **random_state** : *int or RandomState instance, default=None* |
| | When `shuffle` is True, `random_state` affects the ordering of the |