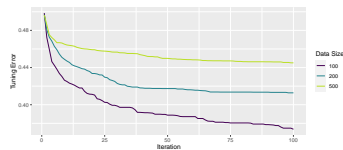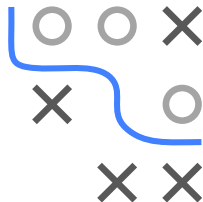# Introduction to Machine Learning
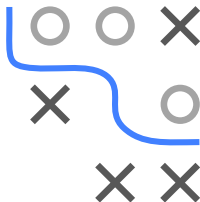
# Nested Resampling
# Motivation



**Learning goals**
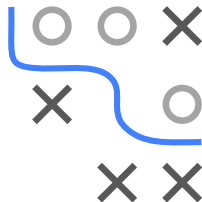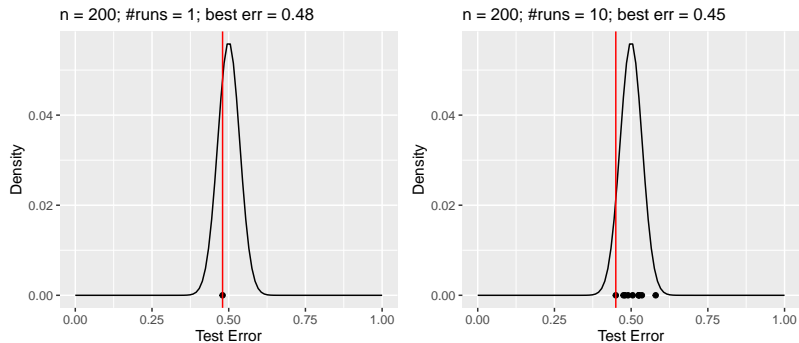
- Understand the problem of overtuning
- Be able to explain the untouched test set principle and how it motivates the idea of nested resampling

# MOTIVATION

Selecting the best model from a set of potential candidates (e.g., different classes of learners, different hyperparameter settings, different feature sets, different preprocessing, ....) is an important part of most machine learning problems.

**Problem**

- We cannot evaluate our finally selected learner on the same resampling splits that we have used to perform model selection for it, e.g., to tune its hyperparameters.
- By repeatedly evaluating the learner on the same test set, or the same CV splits, information about the test set "leaks" into our evaluation.
- Danger of overfitting to the resampling splits / overtuning!
- The final performance estimate will be optimistically biased.
- One could also see this as a problem similar to multiple testing.

# INSTRUCTIVE AND PROBLEMATIC EXAMPLE

- Assume a binary classification problem with equal class sizes.
- Assume a learner with hyperparameter $\lambda$.
- Here, the learner is a (nonsense) feature-independent classifier, where $\lambda$ has no effect. The learner simply predicts random labels with equal probability.
- Of course, its true generalization error is 50%.
- A cross-validation of the learner (with any fixed $\lambda$) will easily show this (given that the partitioned data set for CV is not too small).
- Now let's "tune" it, by trying out 100 different $\lambda$ values.
- We repeat this experiment 50 times and average results.

# INSTRUCTIVE AND PROBLEMATIC EXAMPLE / 2



- Plotted is the best "tuning error" (i.e. the performance of the model with fixed $\lambda$ as evaluated by the cross-validation) after $k$ tuning iterations.
- We have performed the experiment for different sizes of learning data that were cross-validated.

# INSTRUCTIVE AND PROBLEMATIC EXAMPLE



n = 200; #runs = 1; best err = 0.48     n = 200; #runs = 10; best err = 0.45
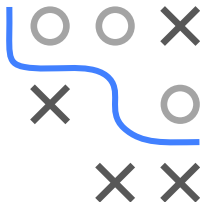
- For 1 experiment, the CV score will be nearly 0.5, as expected
- We basically sample from a (rescaled) binomial distribution when we calculate error rates
- And multiple experiment scores are also nicely arranged around the expected mean 0.5

- But in tuning we take the minimum of those! So we don't really estimate the "average performance" anymore, we get an estimate of "best case" performance instead.
- The more we sample, the more "biased" this value becomes.

# UNTOUCHED TEST SET PRINCIPLE

Countermeasure: simulate what actually happens in model application.

- All parts of the model building (including model selection, preprocessing) should be embedded in the model-finding process **on the training data**.

- The test set should only be touched once, so we have no way of "cheating". The test data set is only used once *after* a model is completely trained, after deciding, for example, on specific hyperparameters.
  Only if we do this are the performance estimates we obtained from the test set **unbiased estimates** of the true performance.

# UNTOUCHED TEST SET PRINCIPLE / 2

- For steps that themselves require resampling (e.g., hyperparameter tuning) this results in **nested resampling**, i.e., resampling strategies for both
    - tuning: an inner resampling loop to find what works best based on training data
    - outer evaluation on data not used for tuning to get honest estimates of the expected performance on new data