

## TK8117 – Week 04 - Topic 01:

# Cross Model Validation

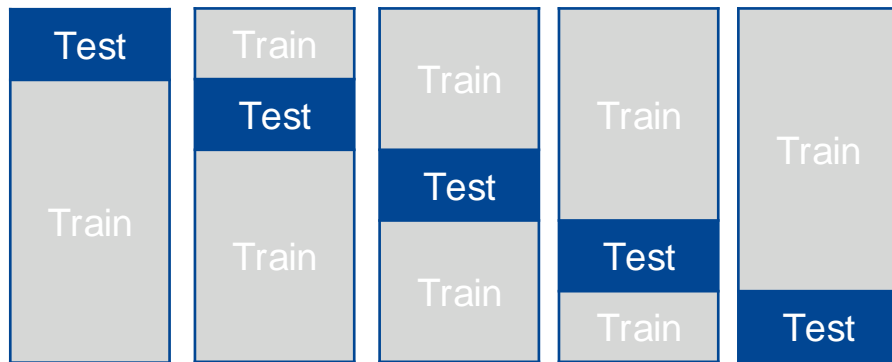
Andreas Gudahl Tufte  
Anette Fagerheim Bjerke  
Aria Alinejad  
Azimil Gani Alam  
Cameron Louis Penne  
Chinmay Anirudha Patwardhan  
Irene Hofmann  
Nicole Quattrini  
Torstein Nordgård-Hansen

# Overview

- Cross-Validation
- Cross *Model* Variation
- Example
- When to use?
- Pros and Cons
- Summary

# Cross-Validation

*Data Partitioning technique into some subsets*



1. Split data into  $k$  folds (e.g.: 5)
2. Train model with training data
3. Measure score with test data
4. Average scores

## Problem:

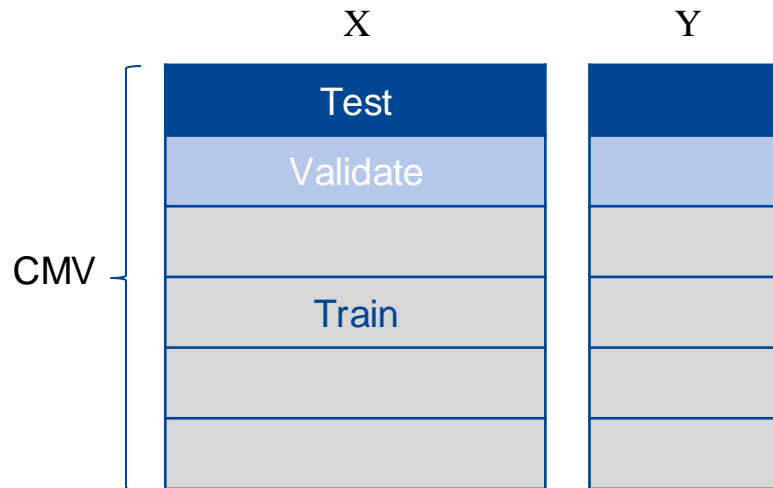
- Risky to find best model by just trying out different variables on the same test data.  
—> high chance that model does not work well on new unseen data

**Solution:** test the model on a completely separate data set.

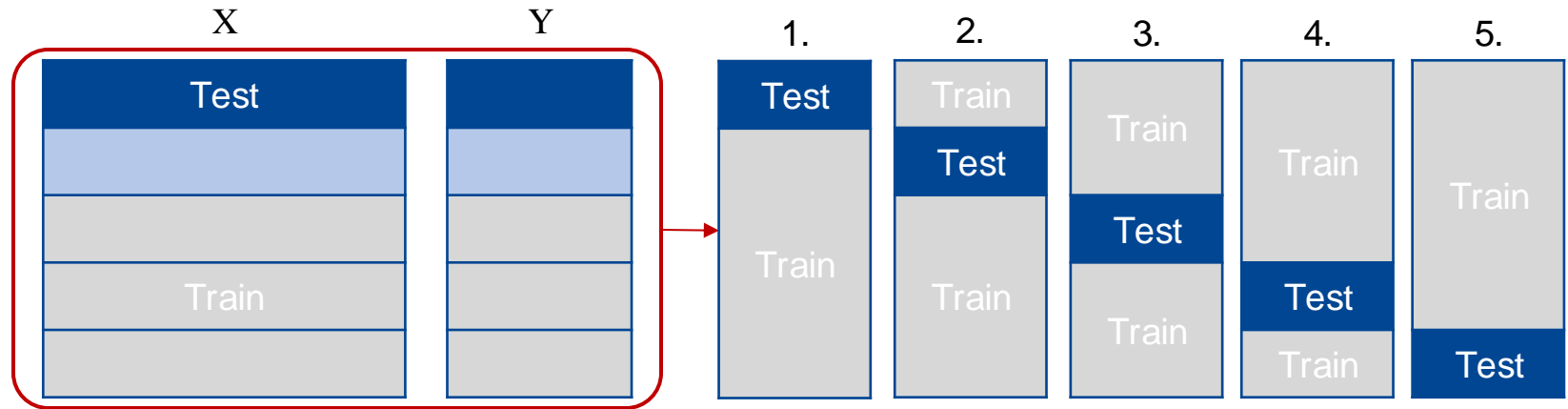
# Cross Model Validation

Also known as:

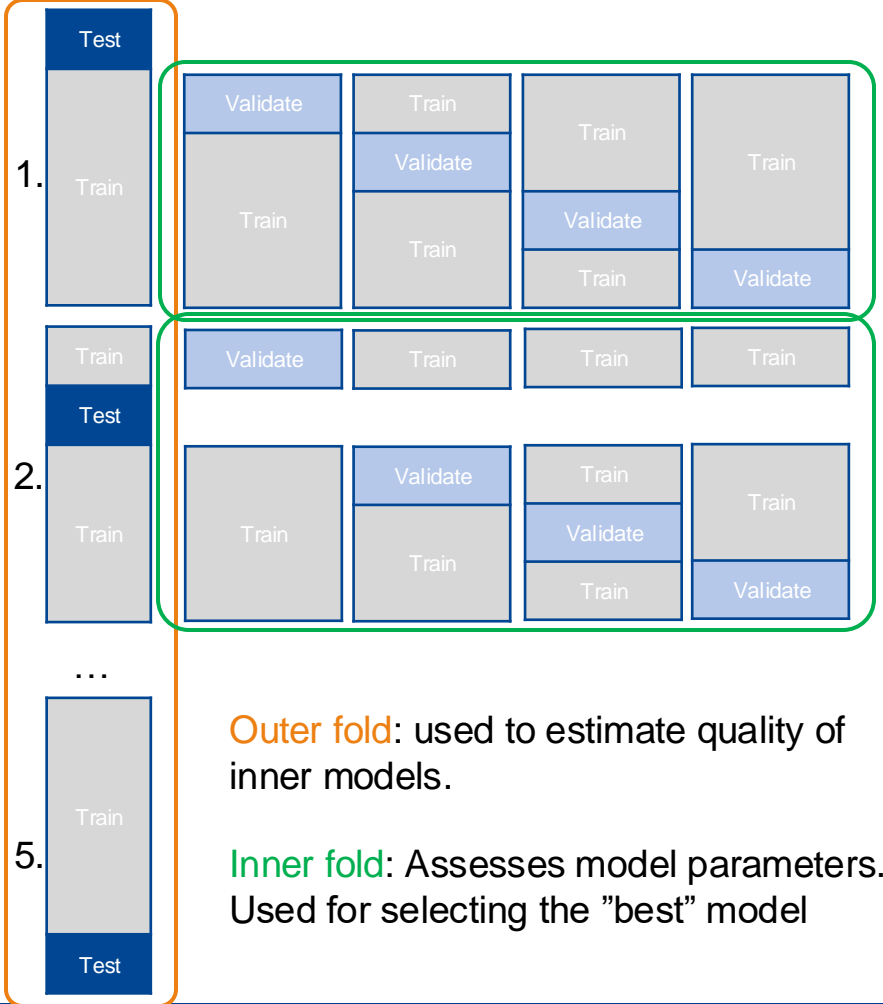
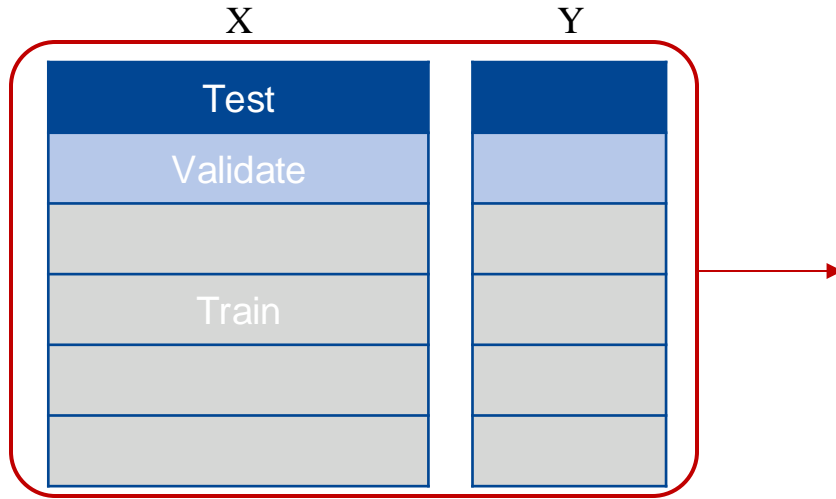
- *Double cross validation*
- *Nested cross validation*
- *Repeated cross validation*



# Cross Model Validation



# Cross Model Validation



**Outer fold:** used to estimate quality of inner models.

**Inner fold:** Assesses model parameters. Used for selecting the "best" model

# Variable Selection

1. Run model with jackknifing
  1. Leaving one sample out
2. Find significant variables using model
  1. PLSR: variable importance in projection (VIP) using weights
3. Repeat 1) and 2) with different sample left out
4. Check frequency of variables found significant during inner fold runs

# Code example

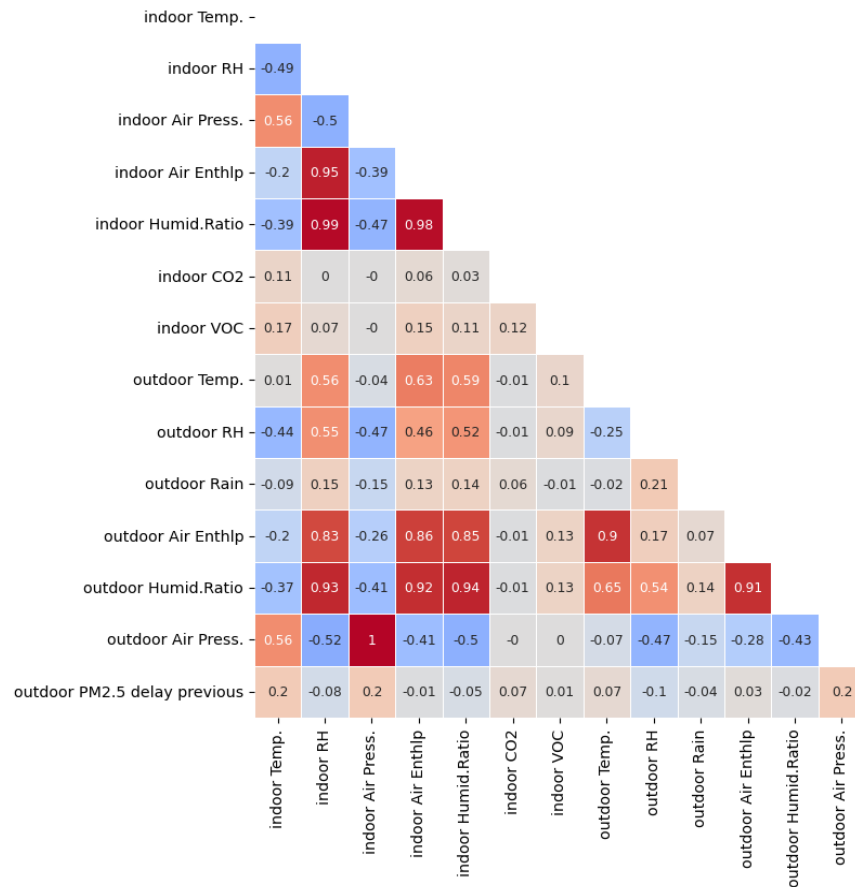


# Data

- Air quality
- Predicting PM1
- 14 variables in dataset

(Unit)	Parameters	measurement location
(ppm)	Carbondioxide (CO2)	indoor
(ppb)	Volatile Organic Compounds (TVOC)	indoor
(ug/m3)	PM1.0	indoor
(ug/m3)	PM2.5	outdoor
(%)	Air Relative Humidity	indoor & outdoor
(deg C)	Air Temperature	indoor & outdoor
(hPa)	Air Pressure	indoor & outdoor
(kJ/kg)	Air Enthalpy	indoor & outdoor
(gram/kg air)	Air Humidity Ratio	indoor & outdoor
(mm)	Rain (Precipitation)	outdoor

Correlation Matrix of Features



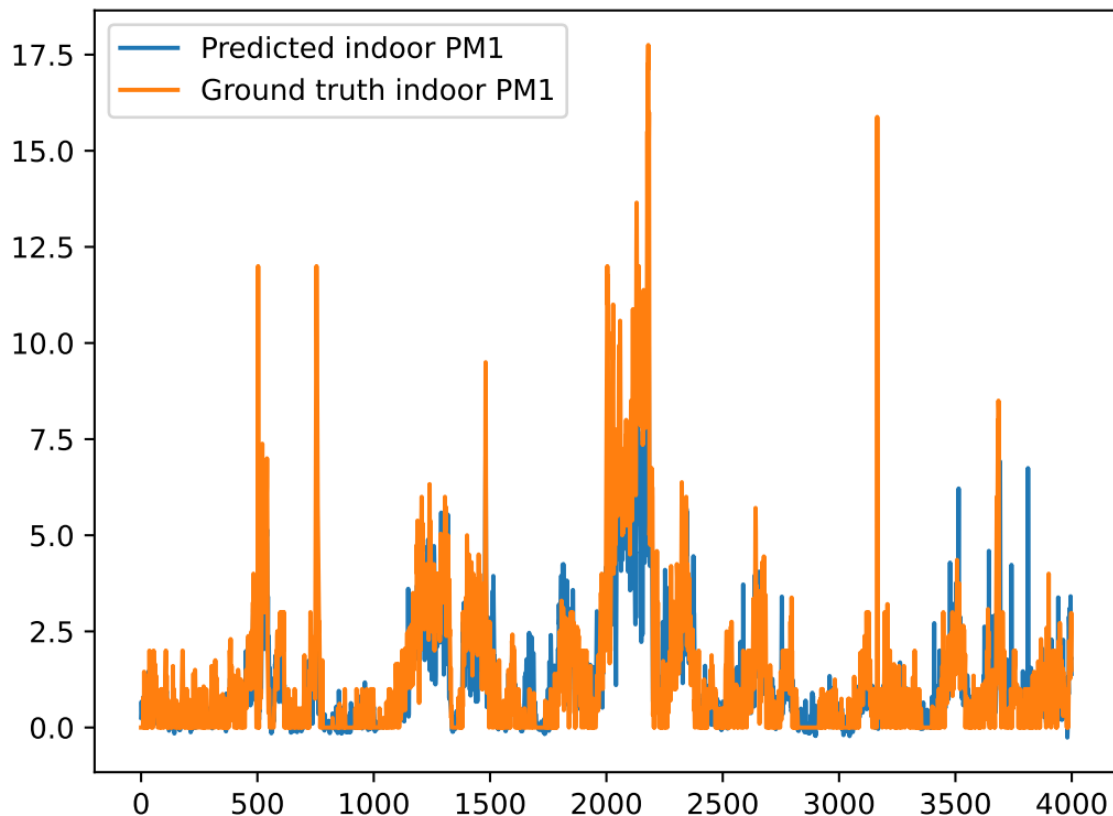
```

55 # 3. Validation and nested cross-validation of the model
56 non_nested_scores = np.zeros(num_trials)
57 nested_scores = np.zeros(num_trials)
58
59 scaler = StandardScaler()
60 scaler.fit(x.values, y=y.values[:, 0])
61 scaled_x = scaler.transform(x.values)
62 scaled_y = y.values[:, 0]
63
64 p_grid = {"C": [1, 10, 100], "gamma": [0.01, 0.1]}
65
66 for i in range(num_trials):
67     # Choose cross-validation techniques for the inner and outer loops
68     inner_cv = KFold(n_splits=4, shuffle=True, random_state=i)
69     outer_cv = KFold(n_splits=4, shuffle=True, random_state=i)
70
71     # Non-nested parameter search and scoring
72     regressor = GridSearchCV(estimator=svm.SVR(), param_grid=p_grid, cv=outer_cv, n_jobs=-1, scoring="r2")
73     regressor.fit(scaled_x, scaled_y)
74     non_nested_scores[i] = regressor.best_score_
75
76     # Nested CV with parameter optimization
77     regressor = GridSearchCV(estimator=svm.SVR(), param_grid=p_grid, cv=inner_cv, n_jobs=-1, scoring="r2")
78     nested_score = cross_val_score(regressor, X=scaled_x, y=scaled_y, cv=outer_cv, n_jobs=-1)
79     nested_scores[i] = nested_score.mean()
80
81 score_difference = non_nested_scores - nested_scores
82

```

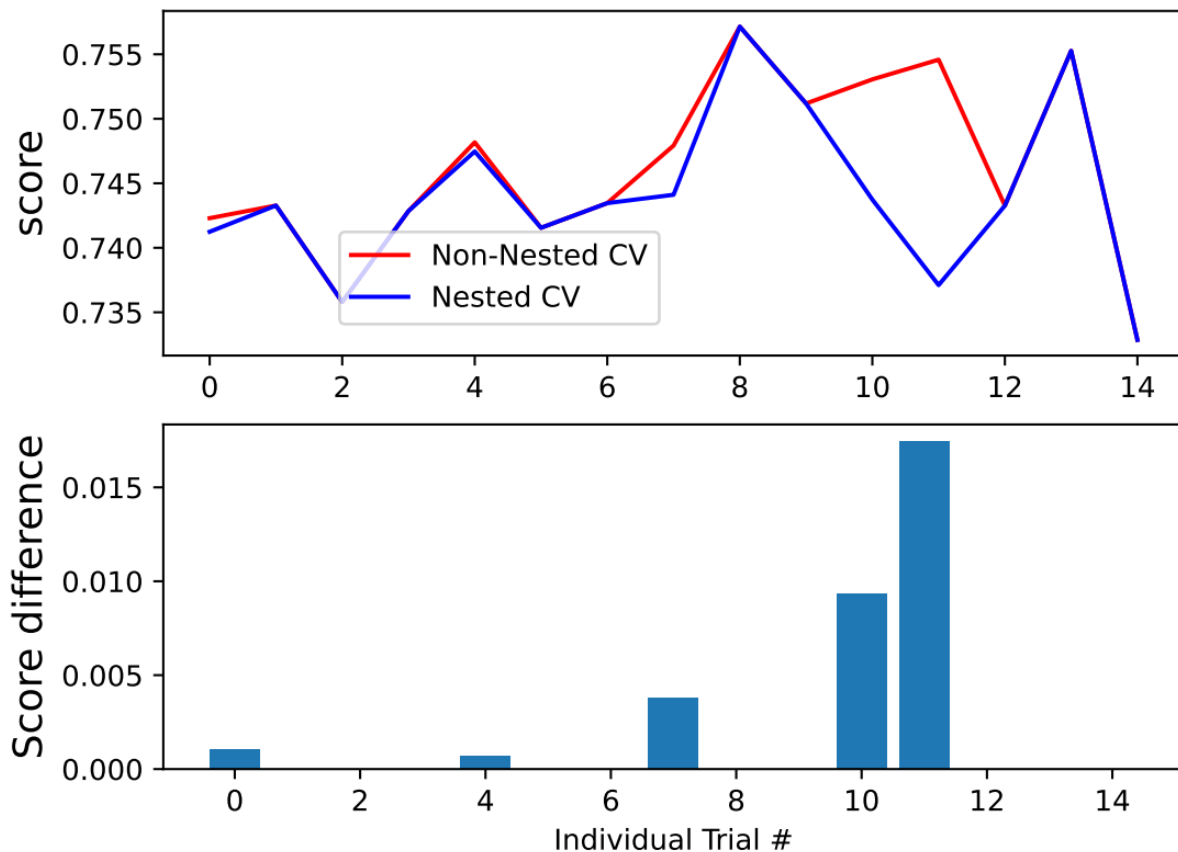
# Predictions

- X-axis: samples
- Y-axis: PM1 value



# Scores

## Non-Nested and Nested Cross Validation on Iris Dataset



# When to use

- Datasets where there are few samples but many variables or hyperparameters
  - Reduce overfitting and false positives
- Identifying subset of variables for model optimization
  - Variable selection

# Pros & Cons

- Prevents overfitting
- Simpler for comparing ML models
- Better hyperparameter tuning than simple cross-validation
- Works well with small data sets
- Slower process due to multiple process inside (special for large datasets),
- Complex setup and interpretation.
- Computationally very resource intensive

# Summary

- **Two-level validation** – involves an outer loop for model evaluation and an inner loop for hyperparameter tuning.
- **Prevents overfitting** – by using separate data for tuning and testing
- **Better Model Performance Estimate** – since the test data in the outer loop is not exposed during the hyperparameter tuning in the inner loop
- **Higher computational cost** – due to multiple model training runs in the inner and outer loops.