# PSY9511: Seminar 4

Model selection, validation and testing
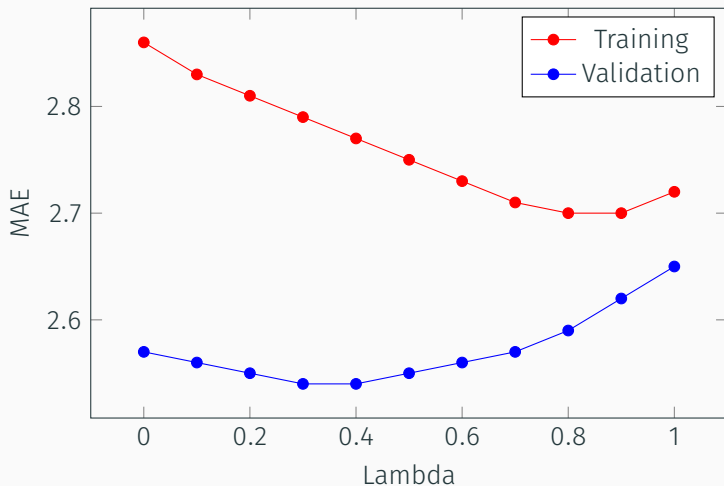
Esten H. Leonardsen

23.09.24

# Outline

1. Assignment 3
2. Loss functions and performance metrics
3. Strategies for model evaluation
    - Training and validation split
    - (Stratification)
    - (Leave-one-out cross-validation)
    - Cross-validation
    - Bootstrap
    - Model comparison
4. Strategies for model selection **and** evaluation
    - Train/validation/test split
    - Nested cross-validation

# Assigment 3

# Assignment 3: Preprocessing

Training

| X |
|---|
| 0.0 |
| 0.2 |
| 0.4 |
| 0.6 |
| 0.8 |
| 1.0 |

Validation

| X |
|---|
| 1.0 |
| 1.2 |
| 1.4 |
| 1.6 |
| 1.8 |
| 2.0 |

Training

| x |
|---|
| 0.0 |
| 0.2 |
| 0.4 |
| 0.6 |
| 0.8 |
| 1.0 |

$\bar{x} = 0.5$

Validation

| x |
|---|
| 1.0 |
| 1.2 |
| 1.4 |
| 1.6 |
| 1.8 |
| 2.0 |

$\bar{x} = 1.5$

Training

| x |
|---|
| 0.0 |
| 0.2 |
| 0.4 |
| 0.6 |
| 0.8 |
| 1.0 |

$\bar{x} = 0.5$

Validation

| x |
|---|
| 1.0 |
| 1.2 |
| 1.4 |
| 1.6 |
| 1.8 |
| 2.0 |

$\bar{x} = 1.5$

In a real-life scenario we don't have a validation set, only individual observations

# Assignment 3: Preprocessing

Training

| X |
|---|
| 0.0 |
| 0.2 |
| 0.4 |
| 0.6 |
| 0.8 |
| 1.0 |

$\bar{x} = 0.5$

Validation

| X |
|---|
| 1.0 |
| 1.2 |
| 1.4 |
| 1.6 |
| 1.8 |
| 2.0 |

$\bar{x} = 1.5$

Training

| x |
|---|
| 0.0 |
| 0.2 |
| 0.4 |
| 0.6 |
| 0.8 |
| 1.0 |

↓

| x |
|---|
| -0.5 |
| -0.3 |
| -0.1 |
| 0.1 |
| 0.3 |
| 0.5 |

Validation

| x |
|---|
| 1.0 |
| 1.2 |
| 1.4 |
| 1.6 |
| 1.8 |
| 2.0 |

↓

| x |
|---|
| -0.5 |
| -0.3 |
| -0.1 |
| 0.1 |
| 0.3 |
| 0.5 |

# The mortal sins of machine learning in neuroscience

1. Reporting inflated modelling performance due to insufficient testing procedures
2. Using accuracy as a model performance measure when it is not appropriate

# The mortal sins of machine learning in neuroscience

1. Reporting inflated modelling performance due to insufficient testing procedures

2. Using accuracy as a model performance measure when it is not appropriate

We believe that we know things that we actually don't (e.g. what brain regions are involved in a mental disorder)

# Loss functions and performance metrics

## Commonalities

- Allows us to evaluate the performance of a model
- Typically on the form $f(y, \hat{y})$

## Loss functions

- Tailored specifically for mathematical optimization of models

## Performance metrics

- Tailored specifically for interpretation of model performance by humans

$$\text{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$
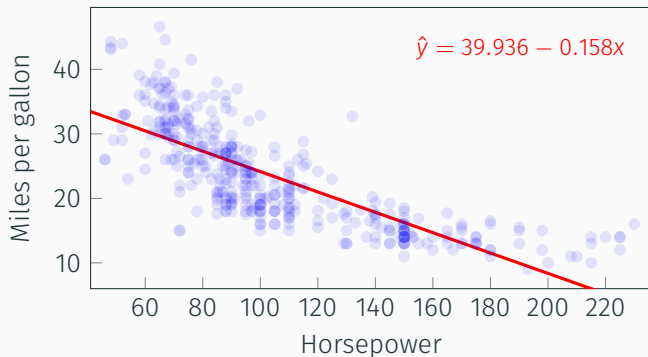
# Loss functions and performance metrics



$$\mathrm{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

$$\hat{y} = 39.936 - 0.158x$$

$$\text{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

$$\text{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

$$\hat{y} = 39.936 - 0.158x$$

$$\mathrm{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

$$= 23.94$$

$$\hat{y} = 39.936 - 0.158x$$

$$\text{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

$$\text{mae}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} |y_i - \hat{y}_i|$$

$$\hat{y} = 39.936 - 0.158x$$
$$\hat{y} = 40.419 - 0.165x$$

$$\text{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

$$\text{mae}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} |y_i - \hat{y}_i|$$

### Loss functions

- Different loss functions measures different properties of the model fit
- Optimizing for them gives different parameter estimates

Tolerance-based accuracy:
A prediction is considered correct if it is within a
predefined margin of error from the true value

$$\text{accuracy}^\star(y, \hat{y}) = \begin{cases} 0 & \text{if } |y - \hat{y}| < \text{tol} \\ 1 & else \end{cases}$$

# Loss functions and performance metrics

| mpg | horsepower |
|-----|------------|
| 22  | 72         |

| mpg | horsepower |
|-----|------------|
| 22  | 72         |

| mpg | horsepower |
|------|------------|
| 22   | 72         |

$$\hat{y} = \beta_0 + \beta_1 \times \text{horsepower}$$

$$\hat{y} = 0 + \beta_1 \times \text{horsepower}$$

$$\hat{y} = 0 + 0.33 \times \text{horsepower}$$

$\hat{y} = 0 + 0.33 \times \text{horsepower}$

### Loss functions

- Different loss functions measures different properties of the model fit
- Optimizing for them gives different parameter estimates
- Must be differentiable to allow for mathematical optimization

$$\text{mse}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

OR

$$\text{mae}(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n} |y_i - \hat{y}_i|$$

$$\frac{1}{n} \sum_{i=0}^{n} (y_i - \hat{y}_i)^2$$

Mean squared error (MSE)

+ Can be used as a loss function

+ Widely used

+ Intuitive

+ Penalizes large errors

? Interpretation

- Depends on scale

$$\sqrt{\frac{1}{n}\sum_{i=0}^{n}(y_i - \hat{y}_i)^2}$$

### Root mean squared error (RMSE)

+ Can be used as a loss function

+ Intuitive

+ Penalizes large errors

+ More interpretable than MSE,
  total loss $\approx$ individual loss

- Depends on scale

$$\frac{1}{n} \sum_{i=0}^{n} |y_i - \hat{y}_i|$$

Mean absolute error (MAE)

- + Can be used as a loss function
- + More interpretable than MSE/RMSE, total loss = average error
- - Feels a bit off
- - Depends on scale

$$\frac{\sum\limits_{i=1}^{n}(y_i-\bar{y})(\hat{y}_i-\bar{\hat{y}})}{\sqrt{\sum\limits_{i=1}^{n}(y_i-\bar{y})^2 \sum\limits_{i=1}^{n}(\hat{y}_i-\bar{\hat{y}})^2}}$$

Pearson correlation coefficient (r)

+ Scale independent
? Captures linear correlation
- Should not be used as a loss function
- Does not care about whether the predictions are close to the true values

$$1 - \frac{\sum\limits_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum\limits_{i=1}^{n}(y_i - \bar{y}_i)^2}$$

### Proportion of variance explained ($r^2$)

+ Scale independent

+ Interpretable

? Captures linear correlation

- Should not be used as a loss function

- Does not care about whether the predictions are close to the true values

Patients

Controls

$y \in \{Patients, Controls\}$

$$y \in \{Patients, Controls\}$$
$$\hat{y} \in \{Patients, Controls\}$$

$$y \in \{0, 1\}$$
$$\hat{y} \in [0, 1]$$

$$y \in \{0, 1\}$$
$$\hat{y} \in \{0, 1\}$$

| TN | FP |
|----|----|
| FN | TP |

Confusion matrix:

Predicted

|  | | 0 | 1 |
|---|---|---|---|
| True | 0 | TN | FP |
| | 1 | FN | TP |

Binary classification metrics:

- Many metrics rely on thresholding the predictions to obtain binary predictions.
- Although not a metric per se, the confusion matrix is a very useful tool to understand model behaviour, and should **always** be looked at (and preferably reported).

$$-(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$$

Logloss

+ Does not rely on thresholding
+ Can be used as a loss function (and very often is)
- Not very interpretable

$$\frac{TP+TN}{TP+TN+FP+FN}$$

Accuracy

+ Interpretable
- **Does not account for imbalanced classes**
- Does not account for different costs of misclassification

Patients ($n = 3$)

Modelling pipeline

| 97 | 0 |
|----|---|
| 3  | 0 |

accuracy=97%

Controls ($n = 97$)

$$\frac{TP}{TP+FN}$$

True positive rate (sensitivity)

+ Interpretable, calculates the proportion of cases that are detected

+ Useful when the cost of false negatives is high (Population-wide screening for severe disease)

$$\frac{TN}{TN+FP}$$

True negative rate (specificity)

+ Interpretable, calculates the proportion of controls that are detected

+ Useful when the cost of false positives is high (Intrusive treatment of rare and mild conditions)

$$\frac{TP}{TP+FP}$$

Positive predictive value (PPV, precision)

+ Interpretable, calculates the proportion of
predicted cases that are actually cases

+ Useful when the cost of false positives is high
(Selection of participants for expensive clinical
trials)

$$\frac{\frac{TP}{TP+FN} + \frac{TN}{TN+FP}}{2}$$

Balanced accuracy

+ Interpretable, behaves similarly to regular accuracy.
+ Takes into account imbalanced classes

Prediction

Prediction

| threshold | TPR | FPR |
|-----------|-----|-----|
|           |     |     |
|           |     |     |
|           |     |     |
|           |     |     |
|           |     |     |

Prediction

| threshold | TPR | FPR |
|-----------|-----|-----|
| 0 | 1 | 1 |
| | | |
| | | |
| | | |
| 1 | 0 | 0 |

Prediction

| threshold | TPR | FPR |
|-----------|-----|-----|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| | | |
| | | |
| 1 | 0 | 0 |

Prediction

| threshold | TPR | FPR |
|-----------|------|-----|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| | | |
| 1 | 0 | 0 |

Prediction

| threshold | TPR | FPR |
| --- | --- | --- |
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| 0.35 | 0.2 | 0.0 |
| 1 | 0 | 0 |

| threshold | TPR | FPR |
|-----------|------|------|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| 0.35 | 0.2 | 0.0 |
| 1 | 0 | 0 |

Receiver operating characteristic curve (ROC)

| threshold | TPR | FPR |
|-----------|------|-----|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| 0.35 | 0.2 | 0.0 |
| 1 | 0 | 0 |

| threshold | TPR | FPR |
|-----------|------|-----|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| 0.35 | 0.2 | 0.0 |
| 1 | 0 | 0 |

| threshold | TPR | FPR |
|-----------|-----|-----|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| 0.35 | 0.2 | 0.0 |
| 1 | 0 | 0 |

| threshold | TPR | FPR |
|-----------|------|------|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| 0.35 | 0.2 | 0.0 |
| 1 | 0 | 0 |

| threshold | TPR | FPR |
|-----------|------|------|
| 0 | 1 | 1 |
| 0.15 | 0.95 | 0.5 |
| 0.25 | 0.5 | 0.2 |
| 0.35 | 0.2 | 0.0 |
| 1 | 0 | 0 |

Area under the receiver operating characteristic curve (AUC/AUROC)

- A performance metric that does not rely on a correct classification threshold
- Measures whether the predictions are ranked correctly (e.g. patients have a higher prediction than controls)
- **Handles class imbalance (relatively well)** and is commonly reported in the literature

**Performance metrics and loss functions measure the performance of a predictive model**

- There is a range of alternatives that can be used, each capturing a different aspect of a model's performance
- It is good practice to report more than one metric
- For regression:
  - MSE is a common loss function with nice mathematical properties.
  - MAE is an intuitive performance metric
- For classification:
  - Log-loss is the most common loss function for probabilistic classifiers
  - AUC is a widely used metric that is easy to interpret, handles class imbalance (to some degree), and is not reliant on the choice of classification threshold

http://localhost:8889/notebooks/notebooks%2FClassification%20metrics.ipynb

# Strategies for model evaluation

UNIVERSITETET
I OSLO

<u>Statistical inference:</u>
Goal: In-sample quantification

<u>Predictive modelling:</u>
Goal: Out-of-sample generalization

How can we test how good our model is on **unseen data** and **be certain that performance holds if we present even more new data**

Dataset

Dataset

Training      Validation

In the validation set approach we split the dataset into two subsets (commonly ~80%/20%), use the first for training the model and the second to test its performance.

+ Accurate estimate of out-of-sample error

+ Simple

- Variable results depending on the exact split

- Only uses a subset of data for training models

- Gives a point estimate of the error, without confidence intervals

<u>Stratification:</u>
Ensuring all folds of the dataset are similar with respect to some given characteristics.

Dataset

```
In[1]:   df = ...
```

Dataset



■ Male
■ Female
Age

```
In[1]:    df = ...
```

Dataset

Training                    Validation
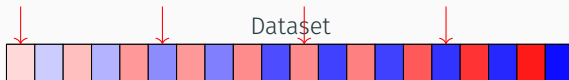
■ Male     ■ Age
■ Female

```
In[1]:    df = ...

          train = df.iloc[:int(len(df) * 0.8)]
          validation = df.iloc[int(len(df) * 0.8):]
```

Dataset

Male
Female
Age

```
In[1]:   df = ...
         df = df.sort_values(['sex', 'age'])
```

# Model evaluation: Stratification



Dataset

```
In[1]:   df = ...
         df = df.sort_values(['sex', 'age'])

         df['fold'] = np.arange(len(df)) % (1 / 0.2)
         train = df[df['fold'] != 0]
         val = df[df['fold'] == 0]
```

Dataset



Training                                    Validation

```
In[1]:   df = ...
         df = df.sort_values(['sex', 'age'])

         df['fold'] = np.arange(len(df)) % (1 / 0.2)
         train = df[df['fold'] != 0]
         val = df[df['fold'] == 0]
```

## Stratification:

Ensuring all folds of the dataset are similar with respect to some given characteristics.

- Helps alleviate the risk of training performance $\gg$ validation performance
- **Always** stratify on target variable first
- Also good idea to stratify on other core characteristics, e.g. sex and age

```
In[1]:   from sklearn.model_selection import  train_test_split
```

```
library(splitstackshape)
stratified(data, columns, split)
```

Dataset

Fits *n* models for *n* datapoints, each time leaving a single datapoint out for testing.

- + Uses all data to train models
- + Not dependent on arbitrary data splits
- + Unbiased (with regards to the full dataset)
- − Computationally expensive
- − Effectively gives a point estimate of the error
- − All models are going to be trained on > 99% overlapping data → highly correlated
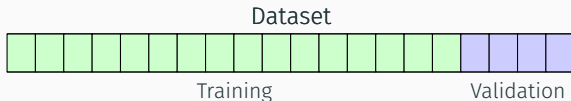
Dataset

Dataset → Error

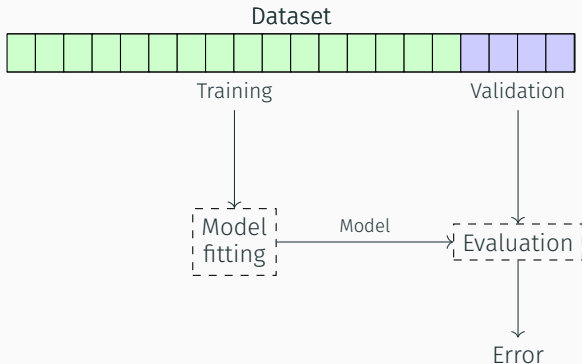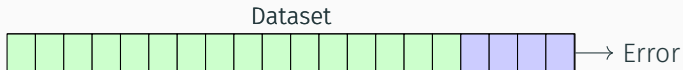Fits $k$ (usually $k \in \{5, 10\}$) models for $n > k$ datapoints, each leaving $n/k$ datapoints for out-of-sample testing.

+ Uses all data to train models
+ Yields multiple estimates of out-of-sample error
- Different choices of $k$ (and exact splits) yields different results
- **No longer a single model from which information (e.g. parameter estimates and p-values) can be derived**

Dataset

Dataset

Dataset

Dataset

Dataset

Dataset

Training

Validation

Dataset

$\longrightarrow$ Error

Fits $b$ models with $m$ datapoints (typically $m < n$), sampled from the original dataset **with replacement**.

+ Uses all data to train models

+ Provides a dense distribution of model performances

+ **Versatile: Can be used for other things, e.g. getting a confidence interval for model parameters**

- Different choices of $b$ (and exact splits) yields different results

Why do we want to evaluate our model?

1. We want to show that our model is better than random guessing
2. We want to show that our model is better than another model

http://localhost:8890/notebooks/notebooks%2FModel%20variability.ipynb

Bad ← Baseline ——— Ours → Good

There is going to be variability in our
model's performance (and possibly the baseline).
**Is our model significantly better?**

Bad ←————————————●————————○————————→ Good
　　　　　　　　　Baseline　　　Ours

Approach 1:
Is the mean of the distribution of performances from our model (with regards to variability that is **unrelated** to efficacy) significantly higher than the point-estimate baseline?

Bad ⟵ ────────────── ⟶ Good

Baseline    Ours

Approach 1:
Is the mean of the distribution of performances from our model (with regards to variability that is **unrelated** to efficacy) significantly higher than the point-estimate baseline?

Approach 1:
Is the mean of the distribution of performances from our model (with regards to variability that is **unrelated** to efficacy) significantly higher than the point-estimate baseline?

Bad ←————————————●————————○————————→ Good
                  Baseline        Ours

<u>Approach 2:</u>
Is the point-estimate performance of our model significantly
higher than the mean of the baseline distribution?

Bad ← ●  ●●●  ●●●● ⬤ → Good

Baseline        Ours

Approach 2:
Is the point-estimate performance of our model significantly
higher than the mean of the baseline distribution?

| Age | Sex | Feature | Outcome |
|-----|--------|---------|---------|
| 25 | Male | 0.53 | 1 |
| 38 | Female | -0.76 | 1 |
| 45 | Male | 0.89 | 1 |
| 33 | Female | -0.21 | 1 |
| 29 | Male | 0.12 | 1 |
| 41 | Female | -0.68 | 0 |
| 56 | Male | 0.45 | 0 |
| 52 | Female | -0.32 | 0 |
| 31 | Male | 0.91 | 0 |
| 48 | Female | -0.15 | 0 |

Modelling pipeline → Error

| Age | Sex | Feature | Outcome |
|-----|--------|---------|---------|
| 25 | Male | 0.53 | 1 |
| 38 | Female | -0.76 | 1 |
| 45 | Male | 0.89 | 1 |
| 33 | Female | -0.21 | 1 |
| 29 | Male | 0.12 | 1 |
| 41 | Female | -0.68 | 0 |
| 56 | Male | 0.45 | 0 |
| 52 | Female | -0.32 | 0 |
| 31 | Male | 0.91 | 0 |
| 48 | Female | -0.15 | 0 |

Modelling pipeline → 0.70

| Age | Sex | Feature | Outcome |
|-----|--------|---------|---------|
| 25 | Male | 0.53 | 1 |
| 38 | Female | -0.76 | 0 |
| 45 | Male | 0.89 | 1 |
| 33 | Female | -0.21 | 0 |
| 29 | Male | 0.12 | 1 |
| 41 | Female | -0.68 | 0 |
| 56 | Male | 0.45 | 1 |
| 52 | Female | -0.32 | 0 |
| 31 | Male | 0.91 | 1 |
| 48 | Female | -0.15 | 0 |

Modelling pipeline → 0.50

| Age | Sex | Feature | Outcome |
|-----|--------|---------|---------|
| 25 | Male | 0.53 | 0 |
| 38 | Female | -0.76 | 0 |
| 45 | Male | 0.89 | 0 |
| 33 | Female | -0.21 | 0 |
| 29 | Male | 0.12 | 0 |
| 41 | Female | -0.68 | 1 |
| 56 | Male | 0.45 | 1 |
| 52 | Female | -0.32 | 1 |
| 31 | Male | 0.91 | 1 |
| 48 | Female | -0.15 | 1 |

Modelling pipeline

0.50
0.45

| Age | Sex | Feature | Outcome |
|-----|--------|---------|---------|
| 25 | Male | 0.53 | 0 |
| 38 | Female | -0.76 | 1 |
| 45 | Male | 0.89 | 0 |
| 33 | Female | -0.21 | 0 |
| 29 | Male | 0.12 | 1 |
| 41 | Female | -0.68 | 1 |
| 56 | Male | 0.45 | 0 |
| 52 | Female | -0.32 | 0 |
| 31 | Male | 0.91 | 1 |
| 48 | Female | -0.15 | 1 |

Modelling pipeline

0.50
0.45
0.55

Bad ←————————— ⬤ ⬤⬤⬤ ⬤⬤⬤⬤ ⬤ —————————→ Good

Baseline            Ours

Approach 2:

Is the point-estimate performance of our model significantly
higher than the mean of the baseline distribution?

Bad ← ... Good

Baseline    Ours

Approach 2:
Is the point-estimate performance of our model significantly
higher than the mean of the baseline distribution?
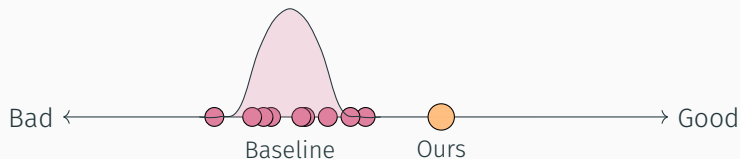
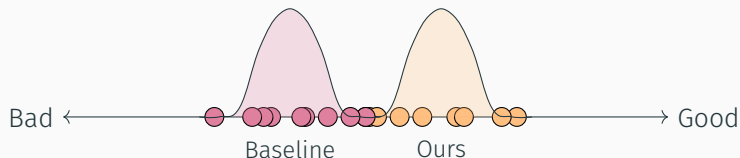Bad ←————————————————————→ Good

Baseline      Ours

Approach 3:

Is the mean of the distribution of performances
from our model significantly higher than the mean
of the distribution of baseline performances?

| Fold | Ours | Baseline |
|------|------|----------|
| 1 | 0.75 | 0.71 |
| 2 | 0.62 | 0.55 |
| 3 | 0.58 | 0.57 |
| 4 | 0.87 | 0.81 |
| 5 | 0.65 | 0.63 |
| 6 | 0.98 | 0.97 |
| 7 | 0.55 | 0.52 |
| 8 | 0.69 | 0.52 |
| 9 | 0.91 | 0.85 |
| 10 | 0.88 | 0.81 |

The small gain of our model will disappear in the noise between the folds using a non-paired statistical test. Use a paired test, e.g. Wilcoxon signed-rank test

- Model evaluation should **always** happen out-of-sample
- If n is big ($\geq$ 10000), a single train/validation split is often sufficient
- For smaller samples, k-fold cross-validation with $5 \leq k \leq 10$ is a good trade-off between bias and variance
- The bootstrap is an effective way of getting confidence intervals for model performance **and parameters**
- Cross-validation (or bootstrapping) will produce a distribution of model performances (although caution the correlation)
- Permutation testing will produce a distribution of baseline performances
- Compare models across folds using Wilcoxon signed-rank test **(ensure the folds are the same!)**

# Model selection and evaluation

- Model evaluation via cross-validation is sufficient if we want to estimate the out-of-sample error of a **known model**.
- Very often we want to know whether a set of predictors are informative for an outcome **given the best possible model**.
- In that case, we have to both **choose the best model**, and **estimate its performance**.
- If we choose the model based on regular cross-validation, the performance estimate will (likely) be inflated

- Model evaluation via cross-validation is sufficient if we want to estimate the out-of-sample error of a **known model**.
- Very often we want to know whether a set of predictors are informative for an outcome **given the best possible model**.
- In that case, we have to both **choose the best model**, and **estimate its performance**.
- If we choose the model based on regular cross-validation, the performance estimate will (likely) be inflated
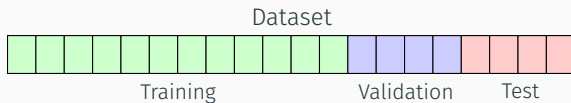- → We need a more advanced strategy
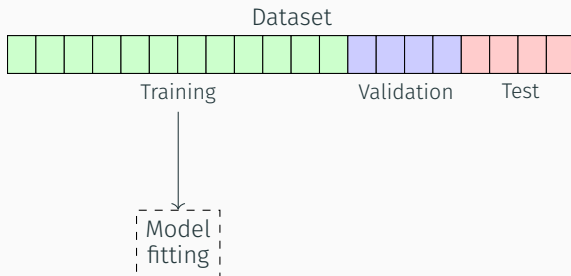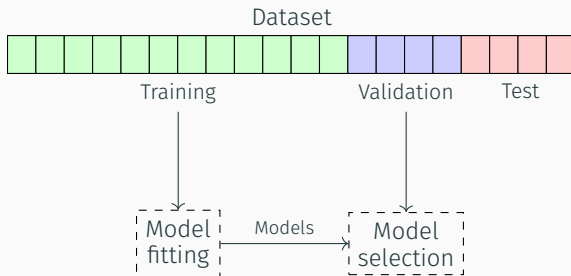
Dataset

Dataset

Training   Validation   Test

Dataset

Training · Validation · Test

Model fitting

Dataset

Dataset

Training and validation                                        Test

Dataset

Training and validation

Test

Error

Error

Error

Error

Error

Average

# Model selection and evaluation: Nested cross-validation

Dataset

Training and validation      Test

Error

Error

Error

Error

Error

Average

Performs $k$ outer cross-validations, each performing $k/k - 1$ inner cross validations, and uses the best models from the inner loop to predict in the outer loop.

+ Uses all data to train models
+ **Unbiased estimate of out-of-sample error**
- **Very** computationally expensive
- We now have either $k$, or $k^2$ models that might behave in different ways

- Whenever a modelling choice is made on the basis of performance in a dataset, **we have to assume the performance achieved by the chosen model is inflated**
- If n is big ($\geq 10000$), a single train/validation/test split is often sufficient
- When possible, use nested cross-validation to select the best model and estimate the out-of-sample error

<span style="color:red">!!!!!</span>

- Whenever a modelling choice is made on the basis of performance in a dataset, **we have to assume the performance achieved by the chosen model is inflated**

- If n is big ($\geq$ 10000), a single train/validation/test split is often sufficient

- When possible, use nested cross-validation to select the best model and estimate the out-of-sample error