

Splitting data

Hold out method

Randomize order of samples before you do any splitting.



→ The model parameters are estimated from the training set.

→ The test set can be used to tune hyperparameters or assess model generalizability (see note after cross-validation below)

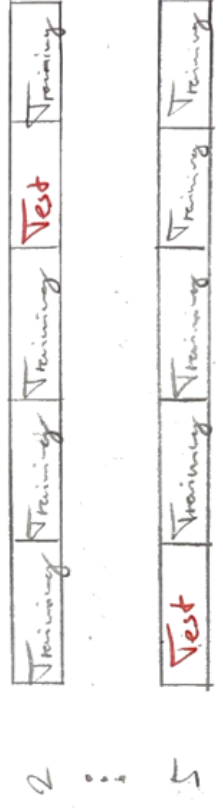
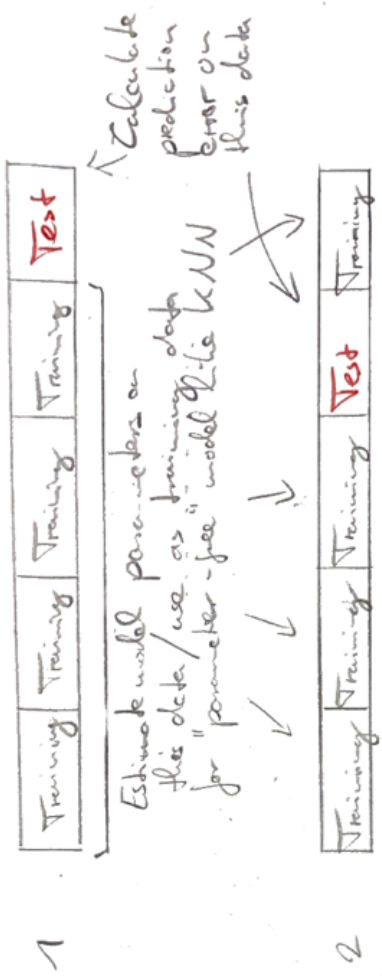
c-fold cross-validation

Randomize order of samples before splitting!



Divide randomized samples into  $c$  equally large subsets, so-called folds

Example: 5-fold cross validation  
Iteration



Notes on tuning parameter / method / model selection

We are looking at a hierarchy of estimation processes:

- Estimation of individual methods / model
- Models with different tuning parameters can be seen as different models (e.g. kNN with  $k=1$  as one model and  $k=5$  another)

Training Data

- Determination which tuning parameters / model / method leads to best performance on test data

Test Data

- In a sense we are performing an estimation process here as well

How do we determine the generalizability of this selected model?

2

In our ideal, data rich world we could create another split.

Training	Validation	Test
----------	------------	------

Estimate individual models

Compare multiple methods

Calculate generalizability of final model

- Prediction error on the test set should be interpreted in the context of the error returned in the validation step.

The test error is much worse than validation then the model doesn't generalize well. However, this doesn't tell you what to change. A more in-depth model/data analysis is necessary.

- In a setting where we do not have enough data to split twice, cross-validation error is usually taken as a surrogate value.

However note that model/method selection can overfit the test set and the prediction error from cross-validation might be lower than what you would get from a separate test set.

## Stratification during data splitting

### Classification example

x x x x  
x x x x  
x x x x  
x x x x  
x x x x  
x x x x

Random splitting of this dataset without paying attention to class labels could easily lead to a training set with only x and no o or a test set without o.

Two class dataset

However, the underlying assumption of the hold out method and cross-validation is that training & test set come from the same distribution.

Having very different class frequencies violates this.

### Stratification as a solution:

class: frequency of x:  $\frac{24}{28} = \frac{6}{7}$

frequency of o:  $\frac{4}{28} = \frac{1}{7}$

Say we want to use 70% of the data for training & 30% for testing and we to maintain class frequencies:

28 data points  $\rightarrow 0.7 \cdot 28 = 19.6$  i.e. 20 training points  
 $\rightarrow 28 - 20 = 8$  test points

Then:  $\frac{6}{7} \cdot 20 \approx 17.1$  i.e. 17 from class x } training  
and 3 from class o }

Summary:

70/30 split for class! Not in the whole dataset  
In test set: 4 from class x  
1 from class o

# 1 Bias - Variance decomposition & tradeoff

For purposes of illustration assume that the response  $y$  fulfills  $E_{P(y|x)}[y] = f(x)$  for some function  $f$ , and  $\text{Var}_{P(y|x)}[y] = \sigma^2$ . (This is like the standard regression model:  $y = f(x) + \epsilon$  with  $E[\epsilon] = 0$ ,  $\text{Var}[\epsilon] = \sigma^2$ )

Then: For the total expected prediction error with squared error loss holds

$$\begin{aligned}
 R &= E_{P(\mathcal{T})} [E_{P(y|x)} [(y - \hat{f}(x))^2]] = \\
 &= E_{P(\mathcal{T}, x)} [E_{P(y|x)} [(y - \hat{f}(x))^2]] = \\
 &= E_{P(\mathcal{T}, x)} [E_{P(y|x)} [(y - E_{P(y|x)}[y]) + (E_{P(y|x)}[y] - \hat{f}(x))^2]] = \\
 &= E_{P(\mathcal{T}, x)} [E_{P(y|x)} [(y - E_{P(y|x)}[y])^2]] + \\
 &\quad + 2 \cdot E_{P(y|x)} [(y - E_{P(y|x)}[y]) (E_{P(y|x)}[y] - \hat{f}(x))] + \\
 &\quad + (E_{P(y|x)}[y] - \hat{f}(x))^2] = \\
 &= \underbrace{E_{P(\mathcal{T}, x)} [E_{P(y|x)} [(y - E_{P(y|x)}[y])^2]]}_{\text{Variance}} + \underbrace{E_{P(\mathcal{T}, x)} [(E_{P(y|x)}[y] - \hat{f}(x))^2]}_{\text{Bias}^2}
 \end{aligned}$$

constant w.r.t.  $f(x)$

1)  $\mathcal{T}$  indep.  $\Rightarrow$   $E_{P(\mathcal{T}, x)} [E_{P(y|x)} [(y - E_{P(y|x)}[y])^2]] = E_{P(y|x)} [(y - E_{P(y|x)}[y])^2] = \text{Var}_{P(y|x)}[y] = \sigma^2$

2)  $\mathcal{T}$  indep.  $\Rightarrow$   $E_{P(\mathcal{T}, x)} [(E_{P(y|x)}[y] - \hat{f}(x))^2] = (E_{P(y|x)}[y] - \hat{f}(x))^2$

$$\begin{aligned}
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(E_{P(y|x)}[y] - \hat{f}(x))^2] = \\
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(y - E_{P(y|x)}[y]) + (E_{P(y|x)}[y] - \hat{f}(x))]^2 - \hat{f}(x)^2 \\
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(y - E_{P(y|x)}[y])^2] + \underbrace{2 \cdot (y - E_{P(y|x)}[y]) (E_{P(y|x)}[y] - \hat{f}(x))}_{\text{Does not depend on } \mathcal{T}} + \\
 &\quad + (E_{P(y|x)}[y] - \hat{f}(x))^2 = \\
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(y - E_{P(y|x)}[y])^2] + \text{Bias}^2
 \end{aligned}$$

$$\begin{aligned}
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(y - E_{P(y|x)}[y])^2] + \text{Bias}^2 \\
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(y - E_{P(y|x)}[y])^2] + \text{Bias}^2 \\
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(y - E_{P(y|x)}[y])^2] + \text{Bias}^2 \\
 &= \sigma^2 + E_{P(\mathcal{T}, x)} [(y - E_{P(y|x)}[y])^2] + \text{Bias}^2
 \end{aligned}$$

