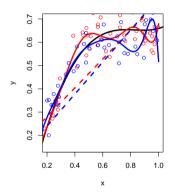
Model Assessment and Hyperparameter Tuning

Johanni Brea

Introduction to Machine Learning

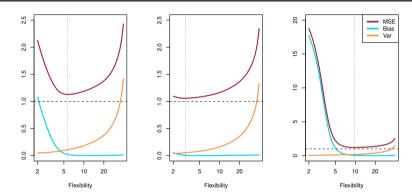


Which Model Is Best?



- Red and blue data points generated with $Y = f(X) + \epsilon$ with $Var(\epsilon) = 0.06^2$ and $f(X) = \sin(2X) + 2(X 0.5)^3 0.5X$ (black line)
- Red/blue lines: fits on red/blue training set
- ► The linear fits (dashed lines) are close to each other (small variance) but far away from f (large bias).
- The polynomial fits (with d = 10) are far from each other (large variance) but close to f (small bias).

Which Model Is Best?



The best model has the smallest test MSE[†].

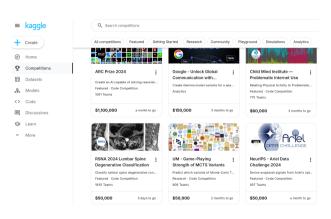
What if we do not know the data generating process?

[†] Here the test MSE is computed exactly, using the data generating process.

Some of the figures in this presentation are taken from "An Introduction to Statistical Learning, with applications in R" (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani



What Would You Do?



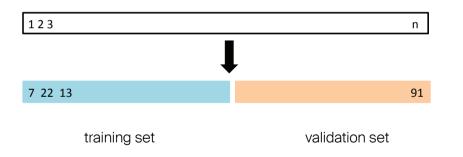
- ▶ You are given training set $\{(x_1, y_1), \dots, (x_n, y_n)\}$ and test inputs $\{x_{n+1}, \dots, x_{n+m}\}$.
- ▶ You test different machine learning methods; each method makes another prediction $\{\hat{y}_{n+1}, \dots, \hat{y}_{n+m}\}$ on the test data.
- ➤ The prediction from which method would you submit to kaggle?

Table of Contents

- 1. Training, Validation and Test Set
- 2. Cross-Validation
- 3. Tuning Models
- 4. A Recipe for Supervised Learning
- 5. The Bootstrap

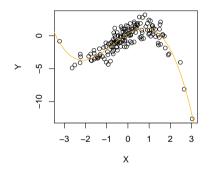
The Validation Set Approach

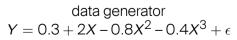
shuffle the data points (each number indicates one row in a data frame) and split into two parts.

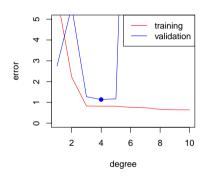




Validation Set Approach Applied to Artificial Data







polynomial fits with different degrees d "optimal" d = 4 (lowest validation error)

Training, Validation and Test Set

- ► **Training Set**: Subset of the full data used to find the parameters.
- ► Validation Set: Held-out subset of the full data used for model selection, i.e. finding the hyper-parameters.
- ▶ Test Set: Held-out subset of the data to estimate the test error of the best model.

Machine Learning Competitions e.g. on kaggle.com

- 1. Start of the competition: Participants obtain a data set, but not the test set.
- 2. Participants split the data set into training and validation sets as they want to fit the parameters and tune the hyper-parameters.
- End of the competition: Organizers evaluate all submitted solutions on the test set.

Recap of Terminology

Assume a simple data generator with Gaussian density $p(y|x) = p(y) = \mathcal{N}(y; \mu, \sigma)$. Assume we found with a machine learning method the predicted mean $\hat{y} = \hat{\mu}$.

- Expected mean squared error: $E[(y-\hat{y})^2] = \int_{-\infty}^{\infty} (y-\hat{y})^2 p(y) dy = \sigma^2 + \mu^2 2\mu \hat{\mu} + \hat{\mu}^2$ Can only be computed, if we know the generator.
- ▶ In practice, we have a dataset $(y_1, ..., y_n)$. We assume it is already shuffled.
 - ► Training loss = approximation of expected mean squared error on training set =

$$\mathcal{L}(\hat{\mu}, \mathcal{D}_{\mathsf{train}}) = \frac{1}{n_{\mathsf{train}}} \sum_{i=1}^{n_{\mathsf{train}}} (y_i - \hat{y})^2 \text{ where } \mathcal{D}_{\mathsf{train}} = (y_1, \dots, y_{n_{\mathsf{train}}}).$$

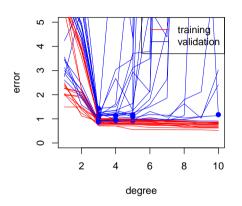
► Test/Validation loss = approximation of expected mean squared error on

test/validation set =
$$\mathcal{L}(\hat{\mu}, \mathcal{D}_{\text{test}}) = \frac{1}{n_{\text{test}}} \sum_{i=n_{\text{train}}+1}^{n} (y_i - \hat{y})^2$$
 where

$$\mathcal{D}_{\mathsf{test}} = (y_{n_{\mathsf{train}}+1}, \dots, y_n)$$
 and $n_{\mathsf{test}} = n - n_{\mathsf{train}}$. Note: $\lim_{n \to \infty} \mathcal{L}(\hat{\mu}, \mathcal{D}_{\mathsf{test}}) = E[(y - \hat{\mu})^2]$



Drawback of Validation Set Approach



The validation estimate of the test error can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the validation set \Rightarrow high variance in model selection.

Quiz

Which of the following statements are correct?

- After finding in a model comparison the best performing model on the validation set, we compute the error on the validation set and the error on the test set.
 - 1. The test error is usually larger than the validation error.
 - 2. Test error and validation error are roughly equal.
 - 3. The test error is usually smaller than the validation error.
- ▶ The error on unseen data tends to be lower for a model trained on all available data compared to a model trained on a training set with 80% of all available data.
- ▶ In a model comparison (e.g. to select hyper-parameters) it is acceptable to fit each model on all available data and compare them on a validation set consisting of 50% of the data.



Table of Contents

- 1. Training, Validation and Test Set
- 2. Cross-Validation
- 3. Tuning Models
- 4. A Recipe for Supervised Learning
- 5. The Bootstrap



Leave-One-Out Cross-Validation (LOOCV)



$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^{n} MSE_{i} \quad MSE_{i} = (y_{i} - \hat{y}_{i})^{2}$$

where \hat{v}_i is the prediction obtained by fitting without (x_i, y_i) . Disadvantage: computational cost of *n* fits (except for linear regression, see section 5.1.2 of textbook).



K-Fold Cross-Validation



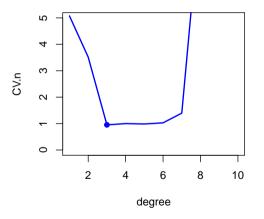
$$CV_{(k)} = \sum_{k=1}^{K} \frac{n_k}{n} MSE_k \quad MSE_k = \frac{1}{n_k} \sum_{i \in C_k} (y_i - \hat{y}_i)^2$$

where \hat{y}_i are predictions obtained by fitting without the data in part C_k .

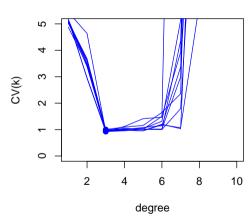


Cross-Validation Applied to the Artificial Data

Leave-one-out Cross Validation

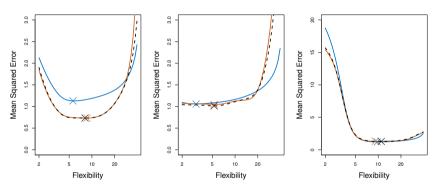


5-Fold Cross Validation





True versus Estimated Test Error



LOOCV (black dashed) and 10-fold CV (orange solid) find almost the same optimal flexibility as the true test error (blue). Crosses indicate the minima of the MSE curves.



Further Considerations

- ▶ The choice of the number of folds *K* is somewhat arbitrary. Typical choices are K = 5 or K = 10.
- ▶ LOOCV has higher computational costs, since n fits are made instead of K (Except for least squares linear or polynomial regression.)
- ▶ To estimate the test error with the validation set approach: fit the winner of model comparison to all data except the test set and evaluate it on the test set.
- ➤ To estimate the test error with cross-validation (nested cross-validation): repeat the approach above for multiple folds.
- ▶ To have the best model for predictions of future data: fit the model on all data vou have.
- ▶ If you do not care about an estimate of the test error, you run cross-validation on the full data, without first splitting off a test set.



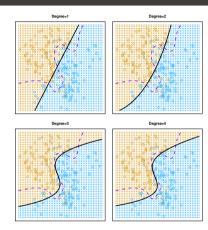
Cross-Validation on Classification Problems

Instead of the log-likelihood, one can also use e.g. the average misclassification rate on the held-out sample for cross-validation in classification problems.

$$\mathsf{CV}_{(n)} = \frac{1}{n} \sum_{i=1}^{n} I(y_i \neq \hat{y}_i)$$

Optimal decision boundary (purple)

Estimated decision boundary (black) for polynomial degrees 1-4.

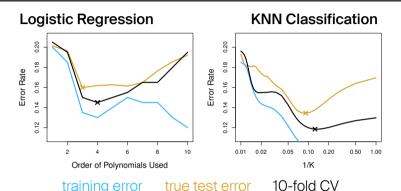


misclassification rates:

degree = 1 : 0.201, degree = 2 : 0.197degree = 3 : 0.160, degree = 4 : 0.162



Cross-Validation on Classification Problems

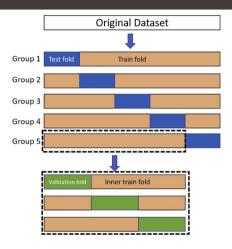


Note: the true test error is known, because the data generator is known in this case.

Side remark: How can it be that the training error goes up with increasing flexibility of the method (order of the polynomial used)? Logistic regression maximizes the likelihood of the parameters β_i given the data \Rightarrow training likelihood is monotonically increasing with order of polynomials, but the training error (misclassification rate) is not necessarily decreasing, because the misclassification rate is a different measure than the log-likelihood.



Nested Cross-Validation



- Outer loop: for each group keep test set for final evaluation.
- Inner loop: find optimal hyper-parameters with standard cross-validation. Note: for each group a different hyper-parameter setting may be optimal.
- Estimate the test error, by computing the average error on the test sets of the outer loop.

figure from https://doi.org/10.1016/j.patter.2021.100329



Quiz

Which of the following statements are correct?

- ► Estimates of the test error with the validation set approach have lower variance than those with LOOCV.
- ▶ In a binary classification task we could use the AUC instead of the error rate to perform cross-validation.



Basic Idea: Split Available Data

Assumption: Hyperparameters fixed Goal: Find best parameters

Use all data to estimate the parameters.

Assumption: Hyperparameters unknown Goal: Find best hyperparameters

Split data into training and validation set(s). Estimate parameters on the training set(s). Estimate test errors for all hyperparameter choices on the validation set(s). Select hyperparameters with lowest test error.

Assumption: Hyperparameters fixed Goal: Estimate test error

Split data into training and test set(s). Estimate parameters on the training set(s). Estimate test error on the test set(s).

Assumption: Hyperparameters unknown Goal: Estimate test error

Split data into training, validation and test set(s). Estimate parameters on the training set(s). Select hyperparameters with lowest test error estimated with the validations set(s). Estimate test error on the test set(s).

Table of Contents

- 1. Training, Validation and Test Set
- 2. Cross-Validation
- 3. Tuning Models
- 4. A Recipe for Supervised Learning
- 5. The Bootstrap



Tuning Models

Hyper-parameter tuning: finding the best model for the given data.

Common recipes:

- ▶ **Grid Search**: Perform cross-validation on a grid of hyper-parameter values. E.g. pick 10 different values of K and pick the best one with cross-validation.
- We will see quite a few hyper-parameter tuning examples in upcoming lectures.
- Use more sophisticated sampling methods for the hyper-parameters to be evaluated, see e.g. https://www.automl.org/or https://github.com/baggepinnen/Hyperopt.jl



Table of Contents

- 1. Training, Validation and Test Set
- 2. Cross-Validation
- 3. Tuning Models
- 4. A Recipe for Supervised Learning
- 5. The Bootstrap



A Recipe for Supervised Learning

- 1. Collect (a lot of) data.
- 2. Look at the raw data; clean it if necessary.
- 3. Select relevant features from the raw data, i.e. choose a suitable representation of the raw data.
- 4. Select a machine learning method.
- 5. Fit the data and tune hyperparameters, e.g. with cross-validation.
- 6. raining loss: high, test loss: high (underfitting?): select a more flexible method.
 - ▶ training loss: low, test loss: high (overfitting?): select a less flexible method.
- 7. Repeat 4-6 until the lowest test loss is found.
- 8. If unhappy with the lowest test loss, repeat 2-7 or collect more data.
- 9. Fit the best model on all available data for best performance on unseen data.



Table of Contents

- 1. Training, Validation and Test Set
- 2. Cross-Validation
- 3. Tuning Models
- 4. A Recipe for Supervised Learning
- 5. The Bootstrap



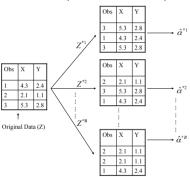
The Bootstrap: a resampling strategy

"Pulling oneself up by one's bootstraps"

19th century saying for impossible tasks



Subsample data with replacement



How can we get new data sets without having new data?

Original data $Z \to B$ bootstrap data sets Z_1^*, \dots, Z_R^*

Some of the figures in this presentation are taken from "An Introduction to Statistical Learning, with applications in R" (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani



Applications of the bootstrap

- ▶ Is it a good idea to obtain a training and a test set with the bootstrap? No, because some samples would appear both in the training and in the test set and therefore the estimate of the test error may be biased (e.g. when overfitting the training set).
- Can we use the bootstrap to estimate uncertainty? Yes. For example, to estimate the variance we should expect, when fitting to different training sets.
 - As a concrete example, say we want to know whether the response (e.g. the wind peak in Luzern) co-varies strongly with a specific feature (e.g. the sunshine duration in Luzern). We could run many linear fits on bootstrapped training sets and look at the distribution of the fitted coefficient (see website).



Suggested Reading

- ▶ 5.1. Cross-Validation
- ▶ 5.2. The Bootstrap

