

A MACHINE LEARNING HANDBOOK

PUBLISHER OF THIS BOOK



L-Università
ta' Malta

Copyright © 2018 ICS5110 APPLIED MACHINE LEARNING class of 2018/9, University of Malta.

JEAN-PAUL EBEJER, DYLAN SEYCHELL, LARA MARIE DEMAJO, **ADD YOUR NAME TO THIS LIST**

Licensed under the Apache License, Version 2.0 (the “License”); you may not use this file except in compliance with the License. You may obtain a copy of the License at <http://www.apache.org/licenses/LICENSE-2.0>. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an “AS IS” BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

First printing, December 2018

Contents

Introduction 5

Cross-Validation 7

Index 13

Introduction

This book explains popular Machine Learning terms. We focus to explain each term comprehensively, through the use of examples and diagrams. The description of each term is written by a student sitting in for ICS5110 APPLIED MACHINE LEARNING¹ at the University of Malta (class 2018/2019). This study-unit is part of the MSc. in AI offered by the Department of Artificial Intelligence, Faculty of ICT.

¹ <https://www.um.edu.mt/courses/studyunit/ICS5110>

Cross-Validation

Cross-validation (CV) is an estimation method used on supervised learning algorithms to assess their ability to predict the output of unseen data [Varma and Simon(2006), Kohavi(1995)]. Supervised learning algorithms are computational tasks like classification or regression, that learn an input-output function based on a set of samples. Such samples are also known as the labeled training data where each example consists of an input vector and its correct output value. After the training phase, a supervised learning algorithm should be able to use the inferred function in order to map new input unseen instances, known as testing data, to their correct output values [Caruana and Niculescu-Mizil(2006)]. When the algorithm incorporates supervised feature selection, cross-validation should always be done external to the selection (feature-selection performed within every CV iteration) so as to ensure the test data remains unseen, reducing bias [Ambroise and McLachlan(2002), Hastie et al.(2001)Hastie, Tibshirani, and Friedman]. Therefore, cross-validation, also known as out-of-sample testing, tests the function's ability to generalize to unseen situations [Varma and Simon(2006), Kohavi(1995)].

Cross-validation has two types of approaches, being i) the exhaustive cross validation approach which divides all the original samples in every possible way, forming training and test sets to train and test the model, and ii) the non-exhaustive cross validation approach which does not consider all the possible ways of splitting the original samples [Arlot et al.(2010)Arlot, Celisse, et al.]. Each of these approaches are further divided into different cross-validation methods, which are explained below.

Exhaustive cross-validation

- Leave- p -out (LpO)

This method takes p samples from the data set as the test set and keeps the remaining as the training set, as shown in Fig. 2a. This is repeated for every combination of test and training set formed from the original data set and the average error is obtained. Therefore, this method trains and tests the algorithm $\binom{n}{p}$ times when the number of samples in the original data set is n , becoming inapplicable when $p > 1$ [Arlot et al.(2010)Arlot, Celisse, et al.].

- Leave-one-out (LOO)

This method is a specific case of the LpO method having $p = 1$. It requires less computation efforts than LpO since the process is only repeated $n_{choose1} = n$ times, however might still be inapplicable for large values of n [Arlot et al.(2010)Arlot, Celisse, et al.].

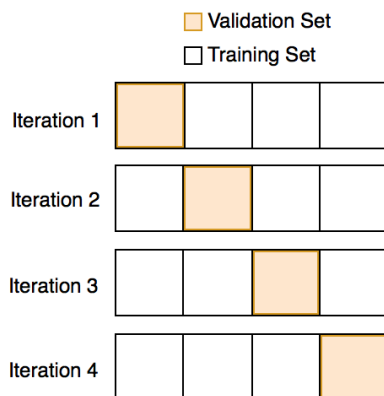
Non-exhaustive cross-validation

- Holdout method

This method randomly splits the original data set into two sets being the training set and the test set. Usually, the test set is smaller than the training set so that the algorithm has more data to train on. This method involves a single run and so must be used carefully to avoid misleading results. It is therefore sometimes not considered a CV method [Kohavi(1995)].

- k -fold

This method randomly splits the original data set into k equally sized subsets, as shown in Fig. 3. The function is then trained and validated k times, each time taking a different subset as the test data and the remaining $(k - 1)$ subsets as the training data, using each of the k subsets as the test set once. The k results are averaged to produce a single estimation. Stratified k -fold cross validation is a refinement of the k -fold method, which splits the original samples into equally sized and distributed subsets, having the same proportions of the different target labels [Kohavi(1995)].



- Repeated random sub-sampling

This method is also known as the Monte Carlo CV. It splits the data set randomly with replacement into training and test subsets using some predefined split percentage, for every run. Therefore, this generates new training and test data for each run but the test data of the different runs might contain repeated samples, unlike that of k -fold [Xu and Liang(2001)].

All of the above cross-validation methods are used to check whether the model has been overfitted or underfitted and hence estimating

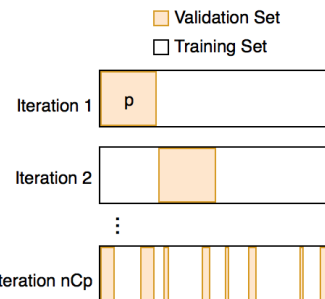


Figure 1: Exhaustive cross-validation methods: Leave-p-Out

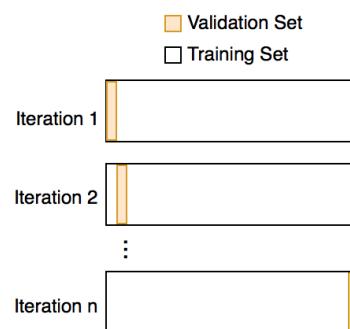


Figure 2: Exhaustive cross-validation methods: Leave-One-Out

Figure 3: k -Fold Cross Validation where $k=4$

the model's ability of fitting to independent data. Such ability is measured using quantitative metrics appropriate for the model and data [Kohavi(1995), Arlot et al.(2010)Arlot, Celisse, et al.]. In the case of classification problems, the misclassification error rate is usually used whilst for regression problems, the mean squared error (MSE) is usually used. MSE is represented by Eq. 1, where n is the total number of test samples, Y_i is the true value of the i^{th} instance and \hat{Y}_i is the predicted value of the i^{th} instance.

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (1)$$

Underfitting is when the model has a low degree (e.g. $y = x$, where the degree is 1) and so is not flexible enough to fit the data making the model have a low variance and high bias [Baumann(2003)], as seen in Fig. 5a. Variance is the model's dependence on the training data and bias is model's assumption about the shape of the data [Arlot et al.(2010)Arlot, Celisse, et al.]. On the other hand, as seen in Fig. 5b, overfitting is when the model has a too high degree (e.g. $y = x^{30}$, where the degree is 30) causing it to exactly fit the data as well as the noise and so lacks the ability to generalize [Baumann(2003)], making the model have a high variance. Cross-validation helps reduce this bias and variance since it uses most of the data for both fitting and testing and so helps the model learn the actual relationship within the data. This makes cross-validation a good technique for models to acquire a good bias-variance tradeoff [Arlot et al.(2010)Arlot, Celisse, et al.].

As stated in [Kohavi(1995)], the LOO method gives a 0% accuracy on the test set when the number of target labels are equal to the number of instances in the dataset. It is shown that the k -fold CV method gives much better results, due to its lower variance, especially when $k = 10, 20$. Furthermore, R. Kohavi et al. state that the best accuracy is achieved when using the stratified cross-validation method, since this has the least bias.

Therefore, let's take an example using the stratified k -fold cross-validation method with $k = 10$. Let's say that we are trying to solve age group classification, using eight non-overlapping age groups being 0-5, 6-10, 11-20, 21-30, 31-40, 41-50, 51-60, and 61+. We are using the FG-NET labelled data set, which contains around 1000 images of individuals aged between 0 and 69. Before we can start training our model (e.g. CNN), we must divide our data set into training and test subsets and this is where cross validation comes in. Therefore, we start by taking the 1000 images of our data set and splitting them according to their target class. Let us assume we have an equal amount of 125 ($1000/8$) images per class². As depicted in Fig. 6, we can now start forming our 10 folds by taking 10% of each age-group bucket, randomly without replacement. Hence, we will end up with 10 subsets of 100 images that are equally distributed along all age-groups. With these subsets, we can estimate our model's accuracy with a lower bias-variance tradeoff. Since we are using

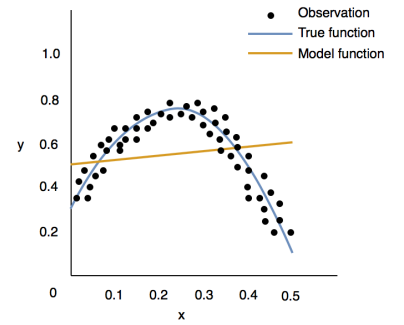


Figure 4: Underfitting

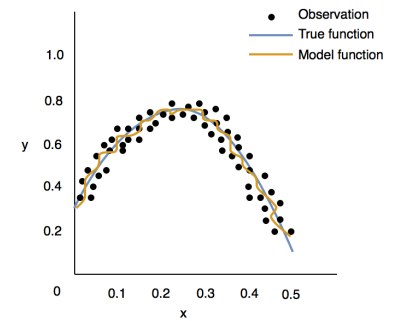


Figure 5: Overfitting

² Down-sampling or up-sampling are common techniques used when there is an unequal amount of samples for the different classes.

10-fold CV, we will train and test our model 10 times. For the first iteration, we shall use subset 1 as the validation set and subsets 2 to 10 as the training set, for the second iteration we use subset 2 as the test set and subsets 1 plus 3 to 10 as our training set, and so on (as shown in Fig. 3). For each iteration we use the misclassification error rate to obtain an accuracy value and we finally average the 10 accuracy rates to obtain the global accuracy of our model when solving age group classification, given the FG-NET data set. Hence, we have now estimated the prediction error of the model and have an idea of how well our model performs in solving such a problem. It is important to note that cross-validation is *just* an estimation method and when using our model in real-life applications we do not apply CV but rather train our model with all the data we have.

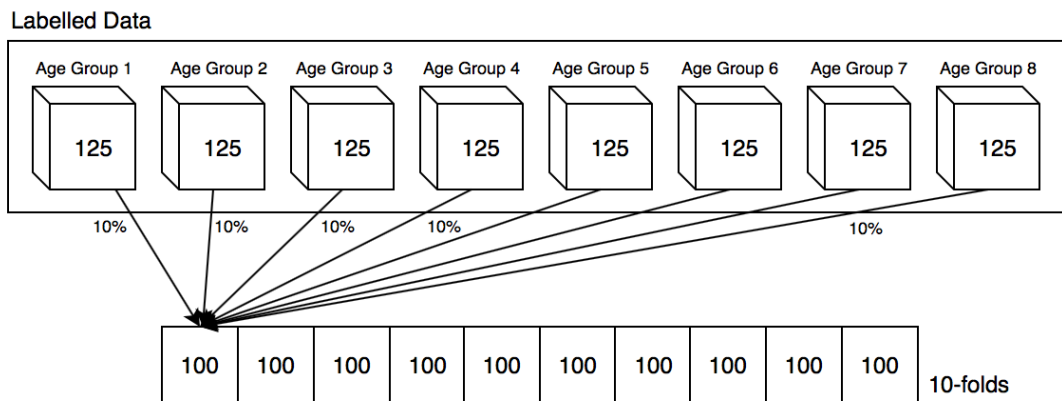


Figure 6: Stratified 10-fold cross-validation on 1000 labelled images of 8 different classes

As concluded by [Varma and Simon(2006)], cross-validation is well implemented when everything is taken place within every CV iteration (including preprocessing, feature-selection, learning new algorithm parameter values, etc.), and the least bias can be achieved when using nested CV methods.

Bibliography

- [Ambroise and McLachlan(2002)] Christophe Ambroise and Geoffrey J McLachlan. Selection bias in gene extraction on the basis of microarray gene-expression data. *Proceedings of the national academy of sciences*, 99(10):6562–6566, 2002.
- [Arlot et al.(2010)Arlot, Celisse, et al.] Sylvain Arlot, Alain Celisse, et al. A survey of cross-validation procedures for model selection. *Statistics surveys*, 4:40–79, 2010.
- [Baumann(2003)] Knut Baumann. Cross-validation as the objective function for variable-selection techniques. *TrAC Trends in Analytical Chemistry*, 22(6):395–406, 2003.
- [Caruana and Niculescu-Mizil(2006)] Rich Caruana and Alexandru Niculescu-Mizil. An empirical comparison of supervised learning algorithms. In *Proceedings of the 23rd international conference on Machine learning*, pages 161–168. ACM, 2006.
- [Hastie et al.(2001)Hastie, Tibshirani, and Friedman] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- [Kohavi(1995)] Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI’95*, pages 1137–1143, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc. ISBN 1-55860-363-8.
- [Varma and Simon(2006)] Sudhir Varma and Richard Simon. Bias in error estimation when using cross-validation for model selection. *BMC bioinformatics*, 7(1):91, 2006.
- [Xu and Liang(2001)] Qing-Song Xu and Yi-Zeng Liang. Monte carlo cross validation. *Chemometrics and Intelligent Laboratory Systems*, 56(1):1–11, 2001.

Index

cross-validation, [7](#)
 holdout, [8](#)
 k-fold, [8](#)
 leave-one-out, [8](#)

 leave-p-out, [7](#)
license, [2](#)
overfitting, [9](#)

repeated random sub-sampling, [8](#)
underfitting, [9](#)