

# *A Machine Learning Study Guide*



# Machine Learning Handbook

***The Definitive Guide***

*ICS5110, class of 2018/9*



**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, DANIEL FARRUGIA, **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

*Confusion Matrix*      7

*Cross-Validation*      11

*Index*      19



# *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 LEARNING<sup>1</sup> 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.

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





# Confusion Matrix

A *Confusion Matrix* (CM), is a contingency table showing how well a model classifies categorical data. By convention (Sammur and Webb, 2017), the CM of an N-class model is an  $N \times N$  matrix indexed by the true class in the row dimension and the predicted class in the column dimension (Table 1).

		Predicted Class	
		<i>spam</i>	$\neg$ <i>spam</i>
True Class	<i>spam</i>	10	1
	$\neg$ <i>spam</i>	2	100

Table 1: CM of a hypothetical binary classifier which predicts whether out-of-sample text objects are spam or not. In this example, 10 spam and 100 non-spam objects are classified correctly, whilst 1 spam and 2 non-spam objects are misclassified.

Even though CMs are commonly used to evaluate binary classifiers, they are not restricted to 2-class models (Martin and Jurafsky, 2018). A CM of a multi-class model would show the number of times the classes were predicted correctly and which classes were confused with each other (Table 2).

	<i>M&amp;M's</i>	<i>Skittles</i>	<i>Smarties</i>
<i>M&amp;M's</i>	34	3	8
<i>Skittles</i>	1	28	5
<i>Smarties</i>	2	4	22

Table 2: CM of a hypothetical sweets classifier. The main diagonal of the CM shows the number of correct predictions, whilst the remaining elements indicate how many sweets were misclassified.

The CM of the model  $h : X \mapsto C$  over the concept  $c : X \mapsto C$  using dataset  $S \subset X$  is formally defined as a matrix  $\Xi$  such that  $\Xi_{c,S}(h)[d_1, d_2] = |S_{h=d_1, c=d_2}|$  (Cichosz, 2014). The CM is constructed by incrementing the element corresponding to the true class *vis-a-vis* the predicted class for each object in the dataset (Algorithm 1).

$\Xi \leftarrow 0$
<b>for</b> $x \in S$ <b>do</b>
$d_1 \leftarrow c(x)$
$d_2 \leftarrow h(x)$
$\Xi_{d_1, d_2} \leftarrow \Xi_{d_1, d_2} + 1$

Algorithm 1: The CM is initialised to the zero matrix, and populated by iterating over all the objects  $x$  with corresponding true class  $d_1$  and predicted class  $d_2$  and incrementing the element  $(d_1, d_2)$  by 1 for each matching outcome.

In binary classification, the CM consists of 2 specially designated classes called the *positive* class and the *negative* class (Saito and Rehmsmeier, 2015). As indicated in Table 3, positive outcomes from the true class which are classified correctly are called *True Positives* (TP), whilst misclassifications are called *False Negatives* (FN). On the

other hand, negative true class outcomes which are classified correctly are called *True Negatives* (TN), and misclassifications are called *False Positives* (FP). In natural sciences, FP are called *Type I Errors* and FN are known as *Type II Errors* (Fielding and Bell, 1997).

	+ve	-ve
+ve	TP	FN
-ve	FP	TN

Table 3: CMs of binary classifiers have positive (+ve) and negative (-ve) classes, and elements called *True Positives* (TP), *False Positives* (FP), *True Negatives* (TN) and *False Negatives* (FN).

The information presented in the CM can be used to evaluate the performance of different binary classifiers (Lu et al., 2004). A number of statistics (Equations 1-7) derived from the CM have been proposed in the literature (Deng et al., 2016) to gain a better understanding of what are the strengths and weaknesses of different classifiers. Caution should be exercised when interpreting metrics (Jeni et al., 2013), since the CM could be misleading if the data is imbalanced and an important subrange of the domain is underrepresented (Raeder et al., 2012). For instance, an albino zebra classifier which always returns negative will achieve high accuracy since albinism is a rare disorder.

These metrics are important in situations in which a particular type of misclassification, i.e. FP or FN, could have worse consequences than the other (Hassanien and Oliva, 2017). For example, FP are more tolerable than FN in classifiers which predict whether a patient has a disease. Both outcomes are undesirable, but in medical applications it is better to err on the side of caution since FN could be fatal.

*Accuracy* (ACC) is the proportion of correct predictions (Equation 1). It is a class-insensitive metric because it can give a high rating to a model which classifies majority class objects correctly but misclassifies interesting minority class objects (Branco et al., 2016). The other metrics should be preferred since they are more class-sensitive and give better indicators when the dataset is imbalanced.

$$ACC = \frac{|TP \cup TN|}{|TP \cup FP \cup TN \cup FN|} \quad (1)$$

*Negative Predictive Value* (NPV) is the ratio of the correct negative predictions from the total negative predictions (Equation 2).

$$NPV = \frac{|TN|}{|TN \cup FN|} \quad (2)$$

*True Negative Rate* (TNR), or *Specificity*, is the ratio of the correct negative predictions from the total true negatives (Equation 3).

$$TNR = \frac{|TN|}{|TN \cup FP|} \quad (3)$$

*True Positive Rate* (TPR), also called *Sensitivity* or *Recall*, is the ratio of the correct positive predictions from the total true positives (Equation 4).

$$TPR = \frac{|TP|}{|TP \cup FN|} \quad (4)$$

Sensitivity and Specificity can be combined into a single metric (Equation 5). These metrics are often used in domains in which minority classes are important (Kuhn and Johnson, 2013). For example, the Sensitivity of a medical classifier (El-Dahshan et al., 2010) measures how many patients with the condition tested positive, and Specificity measures how many did not have the condition and tested negative.

$$\text{Sensitivity} \times \text{Specificity} = \frac{|TP| \times |TN|}{|TP \cup FN| \times |TN \cup FP|} \quad (5)$$

*Positive Predictive Value* (PPV), or *Precision*, is the ratio of the correct positive predictions from the total positive predictions (Equation 6). The difference between accuracy and precision is depicted in Figure 1.

$$\text{PPV} = \frac{|TP|}{|TP \cup FP|} \quad (6)$$

Precision and Recall are borrowed from the discipline of *Information Extraction* (Sokolova and Lapalme, 2009). A composite metric called *F-score*, *F1-score*, or *F-measure* (Equation 7) can be derived by finding their harmonic mean (Kelleher et al., 2015).

$$\text{F-score} = 2 \times \frac{\text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} \quad (7)$$

The complements of ACC, NPV, TNR, TPR and PPV are called, respectively, *Error Rate*, *False Omission Rate*, *false positive rate*, *false negative rate* and *False Discovery Rate*.

The metrics can be adapted for evaluating multi-class models by decomposing an N-class CM into 2-class CMs, and evaluating them individually (Stager et al., 2006). The literature describes two methods for decomposing this kind of CM. In the *1-vs-1* approach, 2-class CMs are constructed for each pairwise class as shown in Table 4.

+ve	-ve
M&M's	{Skittles, Smarties}
Skittles	{M&M's, Smarties}
Smarties	{M&M's, Skittles}

In the *1-vs-rest* approach, 2-class CMs are constructed for each class and the remaining classes combined together as shown in Table 5.

+ve	-ve
M&M's	Skittles $\cup$ Smarties
Skittles	M&M's $\cup$ Smarties
Smarties	Skittles $\cup$ M&M's

Using all metrics could be counterproductive due to information redundancy, but none of the metrics is enough on its own (Ma and Cukic, 2007). For instance, Recall is class-sensitive but it would give a perfect score to an inept model which simply returns the positive class. Thus, the best approach is to evaluate with complementary

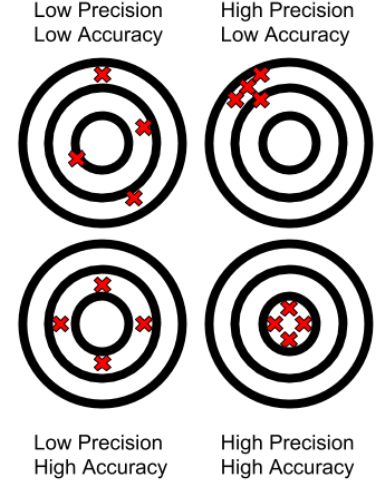


Figure 1: Accuracy vs Precision.

Table 4: 2-class CMs derived from the classes in Table 2. The +ve classes are paired separately with each -ve class.

Table 5: 2-class CMs derived through decomposition of the 3-class CM from Table 2 using the 1-vs-rest approach.

pairs (Gu et al., 2009) such as Sensitivity *vs* Specificity, or Precision *vs* Recall; or a combined measure such as the F-score.

Taking into account the above, CMs are suitable for visualising, evaluating, and comparing the performance of binary or multi-class classifiers. They should be used in conjunction with metrics such as the F-measure to avoid bias, especially if the dataset is unbalanced. For further details on the theoretical aspects of CMs and for practical examples in R refer to (Cichosz, 2014); for examples in Python refer to (Müller et al., 2016).

The following example is motivated by the samples in the *Scikit-Learn* documentation and the work of (Géron, 2017). The models in Figure 2 were trained on the *wines* dataset included with Scikit-Learn.

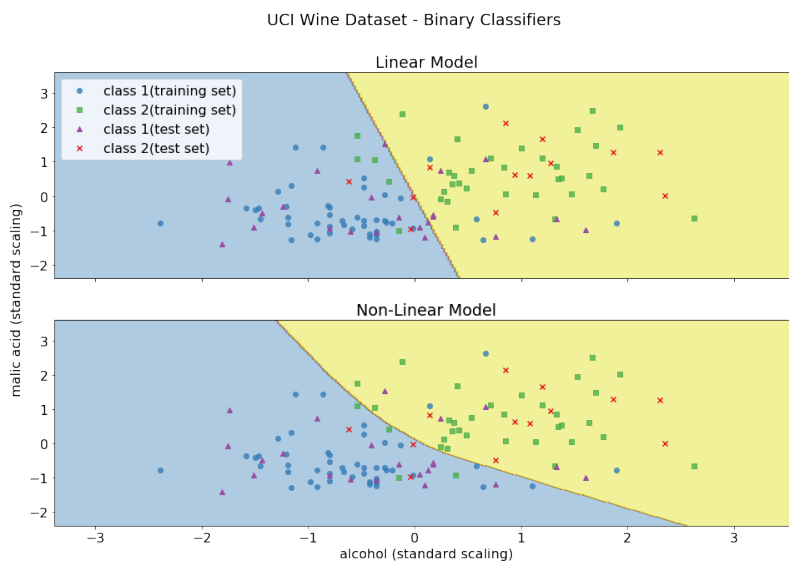


Figure 2: Decision boundary learned by a linear and non-linear binary classifier.

	Linear	Non-Linear
Accuracy	0.72	0.78
Specificity	0.77	0.77
Sensitivity	0.70	0.78
Precision	0.84	0.86
F-score	0.76	0.82

Table 6: Statistics derived from the CMs in Figure 3.

As it can be intuitively deduced from Figure 2, the decision boundary of the non-linear model is a better fit than the linear model. The CMs in Figure 3 also show that non-linear model performs better with a higher TP, and consequently lower TN. The biggest advantage of the non-linear model is the higher Sensitivity resulting in a better F-score.

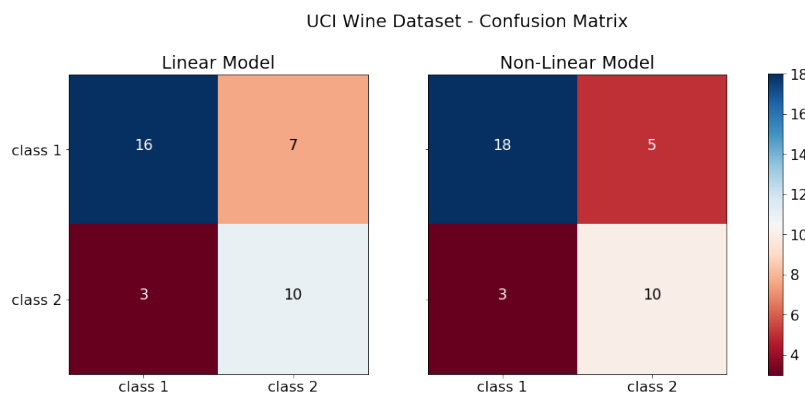


Figure 3: The linear classifier has 16 TP, 10 TN, 7 FN and 3 FP, whilst the non-linear classifier has 18 TP, 10 TN, 5 FN and 3 FP.

# 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). 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).

The above mentioned approaches are further divided into different cross-validation methods, as 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. 5a. 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).

*Leave-one-out (LOO)*

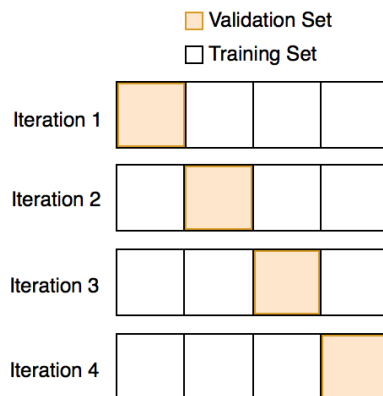
This method is a specific case of the LpO method having  $p = 1$ . It requires less computation efforts than LpO since the process is only repeated  $n_{choose1} = n$  times, however might still be inapplicable for large values of  $n$  (Arlot et al., 2010).

*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. 6. 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

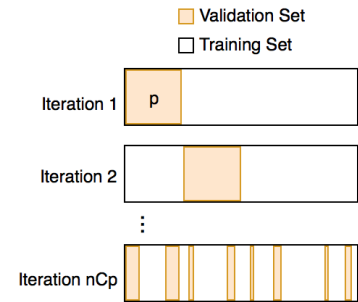


Figure 4: Leave-p-Out Exhaustive Cross Validation

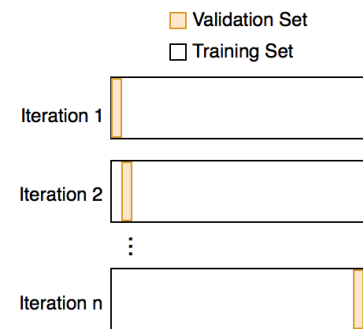


Figure 5: Leave-One-Out Exhaustive Cross Validation

Figure 6:  $k$ -Fold Cross Validation where  $k=4$

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 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). 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. 8, 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.

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (8)$$

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. 8a. 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). On the other hand, as seen in Fig. 8b, 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).

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 class<sup>2</sup>. As depicted in Fig. 9, we can now start forming our 10 folds by taking 10% of each age-group bucket, randomly without replacement. Hence, we will

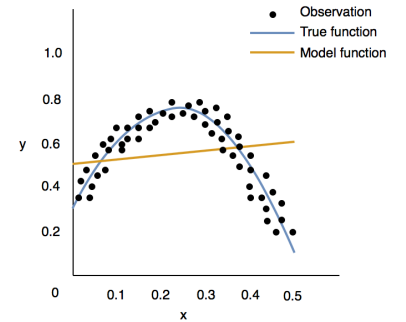


Figure 7: Model Underfitting

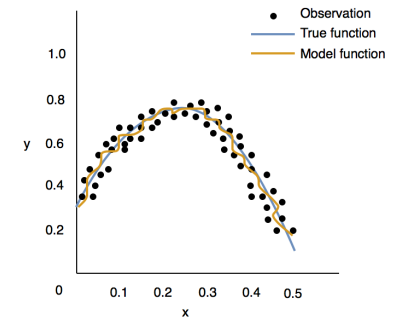


Figure 8: Model Overfitting

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

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 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. 6). 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.

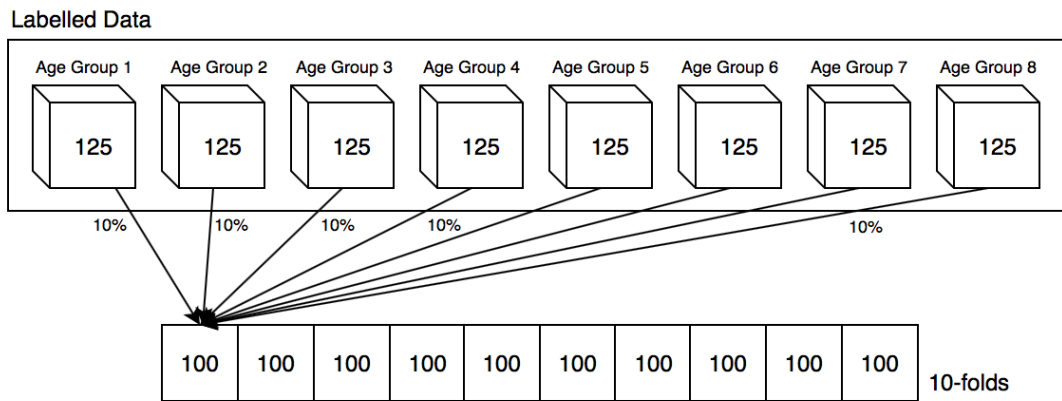


Figure 9: Stratified 10-fold cross-validation on 1000 labelled images of 8 different classes

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

- 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.
- Sylvain Arlot, Alain Celisse, et al. A survey of cross-validation procedures for model selection. *Statistics surveys*, 4:40–79, 2010.
- Knut Baumann. Cross-validation as the objective function for variable-selection techniques. *TrAC Trends in Analytical Chemistry*, 22(6):395–406, 2003.
- Paula Branco, Luís Torgo, and Rita P Ribeiro. A survey of predictive modeling on imbalanced domains. *ACM Computing Surveys (CSUR)*, 49(2):31, 2016.
- 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.
- Pawel Cichosz. *Data mining algorithms: explained using R*. Wiley, 2014.
- Xinyang Deng, Qi Liu, Yong Deng, and Sankaran Mahadevan. An improved method to construct basic probability assignment based on the confusion matrix for classification problem. *Information Sciences*, 340:250–261, 2016.
- El-Sayed Ahmed El-Dahshan, Tamer Hosny, and Abdel-Badeeh M Salem. Hybrid intelligent techniques for mri brain images classification. *Digital Signal Processing*, 20(2):433–441, 2010.
- Alan H Fielding and John F Bell. A review of methods for the assessment of prediction errors in conservation presence/absence models. *Environmental conservation*, 24(1):38–49, 1997.
- Aurélien Géron. *Hands-on machine learning with Scikit-Learn and TensorFlow: concepts, tools, and techniques to build intelligent systems*. O'Reilly, 2017.
- Qiong Gu, Li Zhu, and Zhihua Cai. Evaluation measures of the classification performance of imbalanced data sets. In *International Symposium on Intelligence Computation and Applications*, pages 461–471. Springer, 2009.

- Aboul Ella Hassanien and Diego Alberto Oliva. *Advances in Soft Computing and Machine Learning in Image Processing*, volume 730. Springer, 2017.
- 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.
- László A Jeni, Jeffrey F Cohn, and Fernando De La Torre. Facing imbalanced data—recommendations for the use of performance metrics. In *Affective Computing and Intelligent Interaction (ACII), 2013 Humaine Association Conference on*, pages 245–251. IEEE, 2013.
- John D Kelleher, Brian Mac Namee, and Aoife D’Arcy. *Fundamentals of machine learning for predictive data analytics*. MIT Press, 2015.
- 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.
- Max Kuhn and Kjell Johnson. *Applied predictive modeling*, volume 26. Springer, 2013.
- Zhiyong Lu, Duane Szafron, Russell Greiner, Paul Lu, David S Wishart, Brett Poulin, John Anvik, Cam Macdonell, and Roman Eisner. Predicting subcellular localization of proteins using machine-learned classifiers. *Bioinformatics*, 20(4):547–556, 2004.
- Yan Ma and Bojan Cukic. Adequate and precise evaluation of quality models in software engineering studies. In *Proceedings of the third International workshop on predictor models in software engineering*, page 1. IEEE Computer Society, 2007.
- James H Martin and Daniel Jurafsky. *Speech and language processing*. Pearson, 2018.
- Andreas C Müller, Sarah Guido, et al. *Introduction to machine learning with Python: a guide for data scientists*. O’Reilly, 2016.
- Troy Raeder, George Forman, and Nitesh V Chawla. Learning from imbalanced data: evaluation matters. In *Data mining: Foundations and intelligent paradigms*, pages 315–331. Springer, 2012.
- Takaya Saito and Marc Rehmsmeier. The precision-recall plot is more informative than the roc plot when evaluating binary classifiers on imbalanced datasets. *PloS one*, 10(3):e0118432, 2015.
- Claude Sammut and Geoffrey I Webb. *Encyclopedia of machine learning and data mining*. Springer, 2017.
- Marina Sokolova and Guy Lapalme. A systematic analysis of performance measures for classification tasks. *Information Processing & Management*, 45(4):427–437, 2009.

- Mathias Stager, Paul Lukowicz, and Gerhard Troster. Dealing with class skew in context recognition. In *26th International Conference on Distributed Computing Systems*, pages 58–58. IEEE, 2006.
- Sudhir Varma and Richard Simon. Bias in error estimation when using cross-validation for model selection. *BMC bioinformatics*, 7(1):91, 2006.
- Qing-Song Xu and Yi-Zeng Liang. Monte carlo cross validation. *Chemometrics and Intelligent Laboratory Systems*, 56(1):1–11, 2001.



# *Index*

confusion matrix, [7–10](#)

1-vs-1, [9](#)

1-vs-rest, [9](#)

accuracy, [8](#)

error rate, [9](#)

f-score, [9](#)

f-measure, [9](#)

f1-score, [9](#)

false discovery rate, [9](#)

false negative rate, [9](#)

false omission rate, [9](#)

false positive rate, [9](#)

negative predictive value, [8](#)

positive predictive value, [9](#)

precision, [9](#)

true negative rate, [8](#)

specificity, [8](#)

true positive rate, [8](#)

recall, [8](#)

sensitivity, [8](#)

cross-validation, [11](#)

holdout, [12](#)

k-fold, [12](#)

leave-one-out, [12](#)

leave-p-out, [11](#)

false negatives, [7](#)

false positives, [8](#)

license, [1](#)

overfitting, [13](#)

repeated random sub-sampling, [12](#)

true negatives, [8](#)

true positives, [7](#)

underfitting, [13](#)