# Time for a Change: a Tutorial for Comparing Multiple Classifiers Through Bayesian Analysis

Alessio Benavoli<sup>†</sup> Giorgio Corani<sup>†</sup> Janez Demšar<sup>‡</sup> Marco Zaffalon<sup>†</sup> ALESSIO@IDSIA.CH GIORGIO@IDSIA.CH JANEZ.DEMSAR@FRI.UNI-LJ.SI ZAFFALON@IDSIA.CH

<sup>†</sup> Faculty of Computer and Information Science, University of Ljubljana, Vecna pot 113, SI-1000 Ljubljana, Slovenia <sup>†</sup> Istituto Dalle Molle di Studi sull'Intelligenza Artificiale (IDSIA) Galleria 2, 6928 Manno, Switzerland

#### Editor:

#### Abstract

The machine learning community adopted the use of null hypothesis significance testing (NHST) in order to ensure the statistical validity of results. Many scientific fields however realized the shortcomings of frequentist reasoning and in the most radical cases even banned its use in publications. We should do the same: just as we have embraced the Bayesian paradigm in the development of new machine learning methods, so we should also use it in the analysis of our own results. We argue for abandonment of NHST by exposing its fallacies and, more importantly, offer better—more sound and useful—alternatives for it.

**Keywords:** comparing classifiers, null hypothesis significance testing, pitfalls of p-values, Bayesian hypothesis tests, Bayesian correlated t-test, Bayesian hierarchical correlated t-test, Bayesian signed-rank test

#### 1. Introduction

Progression of Science and of the scientific method go hand in hand. Development of new theories requires—and at the same time facilitates—development of new methods for their validation.

Pioneers of machine learning were playing with ideas: new approaches, such as induction of classification trees, were worthy of publication for the sake of their interestingness. As the field progressed and found more practical uses, variations of similar ideas began emerging, and with that the interest in determining which of them work better in practice. A typical example are the different measures for assessing the quality of attributes; deciding which work better than others required tests on actual, real-world data. Papers thus kept introducing new methods and measured, for instance, classification accuracies to prove their advantages over the existing methods. To ensure the validity of such claims, we adopted—starting with the work of Dietterich (1998) and Salzberg (1997), and later followed by Demšar (2006)—the common statistical methodology used in all scientific areas relying on empirical observations: the null hypothesis significance testing (NHST).

This spread the understanding that the observed results require statistical validation. On the other hand, NHST soon proved inadequate for many reasons (Demšar, 2008). Noteworthy, the American Statistical Association has recently made a statement against p-values (Wasserstein and Lazar, 2016). NHST nowadays it is also falling out of favour in other fields of science (Trafimow and Marks, 2015). We believe that the field of machine learning is ripe for a change as well.

We will spend a whole section demonstrating the many problems of NHST. In a nutshell: it does not answer the question we ask. In a typical scenario, a researcher proposes a new method and desires to prove that it is more accurate than another method on a single data set or on a collection of data sets. She thus runs the competing methods and records their results (classification accuracy or another appropriate score) on one or more data sets, which is followed by NHST. The difference between what the researcher has in mind and what the NHST provides for is evident from the following quote from a recently published paper: "Therefore, at the 90% confidence level, we can conclude that (...) method is able to significantly outperform the other approaches." This is wrong. The stated 90% confidence level is not the probability of one classifier outperforming another. The NHST computes the probability of getting the observed (or a larger) difference between classifiers if the null hypothesis of equivalence was true, which is not the probability of one classifier being more accurate than another, given the observed empirical results. Another common problem is that the claimed statistical significance might have no practical impact. Indeed, the common usage of NHST relies on the wrong assumptions that the p-value is a reasonable proxy for the probability of the null hypothesis and that statistical significance implies practical significance.

As we wrote at the beginning, development of Science not only requires but also facilitates the improvement of scientific methods. Advancement of computational techniques and power reinvigorated the interest for Bayesian statistics. Bayesian modelling is now widely adopted for designing principled algorithms for learning from data (Bishop, 2007; Murphy, 2012). It is time to also switch to Bayesian statistics when it comes to analysis of our own results.

The questions we are actually interested in—e.g., is method A better than B? Based on the experiments, how probably is A better? How high is the probability that A is better by more than 1%?—are questions about posterior probabilities. These are naturally provided by the Bayesian methods (Edwards et al., 1963; Dickey, 1973; Berger and Sellke, 1987). The core of this paper is thus a section that establishes the Bayesian alternatives to frequentist NHST and discusses their inference and results. We eventually describe also the software libraries with the necessary algorithms and give short instructions for their use.

# 2. Frequentist analysis of experimental results

Why do we need to go beyond the frequentist analysis of experimental results? To answer this question, we will focus on a practical case: the comparison of the accuracy of classifiers on different datasets. We initially consider two classifiers: naive Bayes (nbc) and averaged one-dependence estimator (aode). A description of these algorithms with exhaustive references is given in the book by Witten et al. (2011). Assume that our aim is to compare *nbc* versus *aode*. These are the steps we must follow:

- 1. choose a comparison metric;
- 2. select a group of datasets to evaluate the algorithms;
- 3. perform m runs of k-fold cross-validation for each classifier on each dataset.

We have performed these steps in WEKA, choosing accuracy as metric, on a collection of 54 data sets downloaded from the WEKA website<sup>1</sup> and with 10 runs of 10-fold cross-validation. Table 1 reports the accuracies obtained on each dataset by each classifier.

First of all, we aim at knowing which is the best classifier for each dataset. The answer to this question is probabilistic, since on each data set we have only estimates of the performance of each classifier.

Datasets		10 run	s of 10-	fold cros	ss-vali	dation	
anneal	94.44	98.89	94.44	98.89		94.38	97.75
anneal	96.67	100.0	96.67	100.0		96.63	97.75
audiology	73.91	69.56	73.91	60.87		72.73	59.09
audiology	73.91	69.56	78.26	60.87		72.73	59.09
breast-cancer	90.32	90.32	87.1	86.67		86.67	90.0
breast-cancer	87.1	87.1	87.1	86.67		83.33	86.67
$\mathrm{cmc}$	51.35	50.68	54.73	59.18		50.34	48.3
$\mathrm{cmc}$	52.7	50.68	52.7	55.1		52.38	48.98
wine	100.0	95.71	97.14	94.29		97.14	97.1
wine	100.0	95.71	97.14	92.86		97.14	97.1
yeast	57.72	55.03	59.06	58.39		55.4	55.4
yeast	57.05	55.03	59.06	58.39		54.05	55.4
ZOO	81.82	100.0	100.0	90.0		90.0	100.0
ZOO	90.91	100.0	100.0	90.0		90.0	100.0

Table 1: Accuracies for 10 runs of 10-fold cross-validation on 54 UCI datasets for *nbc* (blue row) versus *aode* (white rows).

.

Since during cross-validation we have provided both classifiers with the same training and test sets, we can compare the classifiers by considering the difference in accuracies on each test set. This yields the vector of differences of accuracies  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$ , where n = 100 (10 runs of 10-fold cross-validation). We can compute the mean of the vector of differences  $\mathbf{x}$ , i.e., the mean difference of accuracy between the two classifiers, and statistically evaluate whether the mean difference is significantly different from zero. In frequentist analysis, we would perform a NHST using the t-test. The problem is that the t-test assumes the observations to be independent. However, the differences of accuracies

<sup>1.</sup> See http://www.cs.waikato.ac.nz/ml/weka/.

are not independent of each other because of the overlapping training sets used in cross-validation. Thus the usual t-test is not calibrated when applied to the analysis of cross-validation results: when sampling the data under the null hypothesis its rate of Type I errors is much larger than  $\alpha$  (Dietterich, 1998). Moreover, the correlation cannot be estimated from data; Nadeau and Bengio (2003) have proven that there is no unbiased estimator of the correlation of the results obtained on the different folds. Introducing some approximations, they have proposed a heuristic to choose the correlation parameter:  $\rho = \frac{n_{te}}{n_{tot}}$ , where  $n_{te}$ ,  $n_{tr}$  and  $n_{tot} = n_{te} + n_{tr}$  respectively denote the size of the training set, of the test set and of the whole available data set.<sup>2</sup>

#### Frequentist correlated t-test

The correlated t-test is based on the modified Student's t-statistic:

$$t(\boldsymbol{x},\mu) = \frac{\overline{x} - \mu}{\sqrt{\hat{\sigma}^2(\frac{1}{n} + \frac{\rho}{1-\rho})}} = \frac{\overline{x} - \mu}{\sqrt{\hat{\sigma}^2(\frac{1}{n} + \frac{n_{te}}{n_{tr}})}},$$
(1)

where  $\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$  and  $\hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \overline{x})^2}$  are the sample mean and sample standard deviation of the data x,  $\rho$  is the correlation between the observations and  $\mu$  is the value of the mean we aim at testing. The statistic follows a Student distribution with n-1 degrees of freedom:

$$St\left(\bar{x}; n-1, \mu, \left(\frac{1}{n} + \frac{\rho}{1-\rho}\right)\hat{\sigma}^2\right).$$
 (2)

For  $\rho=0$ , we obtain the traditional t-test. For  $\rho=\frac{n_{te}}{n_{tot}}$ , we obtain the correlated t-test proposed by Nadeau and Bengio (2003) to account for the correlation due to the overlapping training sets. Usually the test is run in a two-sided fashion. Its hypotheses are:  $H_0: \mu=0$ ;  $H_1: \mu\neq 0$ . The p-value of the statistic under the null hypotheses is:

$$p = 2 \cdot (1 - \mathcal{T}_{n-1}(|t(x,0)|)), \tag{3}$$

where  $\mathcal{T}_{n-1}(|t(\boldsymbol{x},0)|)$  denotes the cumulative distribution of the standardized Student distribution with n-1 degrees of freedom in  $|t(\boldsymbol{x},\mu)|$  for  $\mu=0$ . For instance, for the first data set in Table 1 we have that  $\bar{x}=-0.0194$ ,  $\hat{\sigma}=0.01583$ ,  $\rho=1/10$ , n=100 and so  $t(\bar{x},0)=-3.52$ . Hence, the two-sided p-value is  $p=2\cdot(1-\mathcal{T}_{n-1}(|t(\boldsymbol{x},0)|))=0.00065\approx0.001$ . Sometimes the directional one-sided test is performed. If the alternative hypothesis is the positive one, the hypotheses of the one-sided test are:  $H_0: \mu \leq 0$ ;  $H_1: \mu > 0$ . The p-value is  $p=1-\mathcal{T}_{n-1}(t(\boldsymbol{x},0))$ .

Table 2 reports the two-sided p-values for each comparison on each dataset obtained via the correlated t-test. The common practice in NHST is to declare all the comparisons such that  $p \leq 0.05$  as significant, i.e., the accuracy of the two classifiers is significantly different on that dataset. Conversely, all the comparisons with p > 0.05 are declared not significant.

<sup>2.</sup> Nadeau and Bengio (2003) considered the case in which random training and test sets are drawn from the original data set. This is slightly different from k-fold cross-validation, in which the folds are designed not to overlap. However the correlation heuristic by Nadeau and Bengio (2003) has since become commonly used to analyse the cross-validation results (Bouckaert, 2003).

Note that these significance tests can, under the NHST paradigm, only be considered in isolation, while combined they require either an omnibus test like ANOVA or corrections for multiple comparisons.

Dataset	p-value	Dataset	p-value	Dataset	p-value
anneal	0.001	audiology	0.622	breast-cancer	0.598
cmc	0.338	contact-lenses	0.643	$\operatorname{credit}$	0.479
german-credit	0.171	pima-diabetes	0.781	ecoli	0.001
eucalyptus	0.258	$\operatorname{glass}$	0.162	grub-damage	0.090
haberman	0.671	hayes-roth	1.000	cleeland-14	0.525
hungarian-14	0.878	hepatitis	0.048	hypothyroid	0.287
ionosphere	0.684	iris	0.000	kr-s-kp	0.646
labor	1.000	lier-disorders	0.270	lymphography	0.018
monks1	0.000	monks3	0.220	monks	0.000
mushroom	0.000	nursery	0.000	optdigits	0.000
page	0.687	pasture	0.000	pendigits	0.452
postoperatie	0.582	primary-tumor	0.492	segment	0.000
solar-flare-C	0.035	solar-flare-m	0.596	solar-flare-X	0.004
sonar	0.777	soybean	0.049	spambase	0.000
spect-reordered	0.198	splice	0.004	squash-stored	0.940
squash-unstored	0.304	tae	0.684	credit	0.000
owel	0.000	waveform	0.417	white-clover	0.463
wine	0.671	yeast	0.576	ZOO	0.435

Table 2: Two sided p-values for each dataset. The difference is significant (p < 0.05) in 19 out of 54 comparisons.

# 2.1 NHST: the pitfalls of black and white thinking

Despite being criticized from its inception, NHST is still considered necessary for publication, as  $p \leq 0.05$  is trusted as an objective proof of the method's quality. One of the key problems of decisions based on  $p \leq 0.05$  is that it leads to "black and white thinking", which ignores the fact that (i) a statistically significant difference is completely different from a practically significant difference (Berger and Sellke, 1987); (ii) two methods that are not statistically significantly different are not necessarily equivalent. The NHST and this p-value-related "black and white thinking" do not allow for making informed decisions. Hereafter, we list the limits of NHST in order of severity using, as a working example, the assessment of the performance of classifiers.

NHST does not estimate probabilities of hypotheses. What is the probability that the performance of two classifiers is different (or equal)? This is the question we are asking when we compare two classifiers; and NHST cannot answer it.

In fact, the p-value represents the probability of getting the observed (or larger) differences assuming that the performance of the classifiers is equal  $(H_0)$ . Formally, p = p(t(x))

 $\tau|H_0$ ), where  $t(\boldsymbol{x})$  is the statistic computed from the data  $\boldsymbol{x}$ , and  $\tau$  is the critical value corresponding to the test and the selected  $\alpha$ . This is not the probability of the hypothesis,  $p(H_0|\boldsymbol{x})$ , in which we are interested.

Yet, researchers want to know the probability of the null and the alternative hypotheses on the basis of the observed data, rather than the probability of the data assuming the null hypothesis to be true. Sentences like "at the 95% confidence level, we can conclude that (...)", are formally correct, but they seem to imply that 1 - p is the probability of the alternative hypothesis, while in fact  $1 - p = 1 - p(t(\mathbf{x}) > \tau | H_0) = p(t(\mathbf{x}) < \tau | H_0)$ , which is not the same as  $p(H_1|\mathbf{x})$ . This is summed up in Table 3.

# what we compute what we would like to know $\frac{p(t(\boldsymbol{x}) > \tau | H_0)}{1 - p(t(\boldsymbol{x}) > \tau | H_0) = p(t(\boldsymbol{x}) < \tau | H_0)} \qquad \frac{p(H_0 | \boldsymbol{x})}{1 - p(H_0 | \boldsymbol{x}) = p(H_1 | \boldsymbol{x})}$

Table 3: Difference between the probabilities of interest for the analyst and the probabilities computed by the frequentist test.

Point-wise null hypotheses are practically always false. The difference between two classifiers can be very small; however there are no two classifiers whose accuracies are perfectly equivalent.

By using a NHST, the null hypothesis is that the classifiers are equal. However, the null hypothesis is practically always false! By rejecting the null hypothesis NHST indicates that the null hypothesis is unlikely; but this is known even before running the experiment. This problem of the NHST has been pointed out in many different scientific domains (Lecoutre and Poitevineau, 2014, Sec 4.1.2.2). A consequence is that, since the null hypothesis is always false, by adding enough data points it is possible to claim significance even when the effect size is trivial. This is because the p-value is affected both by the sample size and the effect size, as discussed in the next section. Quoting Kruschke and Liddell (2015): "null hypotheses are straw men that can virtually always be rejected with enough data."

The p-value does not separate between the effect size and the sample size. The usual explanation for this phenomenon is that if the effect size  $H_0$  is small, more data is needed to demonstrate the difference. Enough data can confirm arbitrarily small effects. Since the sample size is manipulated by the researcher and the null hypothesis is always wrong, the researcher can reject it by testing the classifiers on enough data. On the contrary, conceivable differences may fail to yield small p-values if there are not enough suitable data for testing the method (e.g., not enough datasets). Even if we pay attention not to confuse the p-value with the probability of the null hypothesis, the p-value is intuitively understood as the indicator of the effect size. In practice, it is the function of effect size and sample size: same p-values do not imply same effect sizes.

Figure 1 reports the *density plots* of the differences of accuracy between nbc and aode on the dataset *hepatitis* in two cases: (1) considering only 15 of the 100 accuracies in Table 1 (left); (2) considering all the 100 accuracies (right). The *two orange vertical lines* 

define the region in which the differences of accuracy is less than 1%—the meaning of these lines will be clarified in the next sections.

The p-value is 0.077 in the first case and so the null hypothesis cannot be rejected. The p-value becomes 0.048 in the second case and so the null hypothesis can be rejected. This demonstrates how adding data leads to rejection of the null hypothesis although the difference between the two classifiers is very small in this dataset (all the mass is inside the two orange vertical lines). Practical significance can be equated with the effect size, which is what the researcher is interested in. Statistical significance—the p-value—is not a measure of practical significance, as shown in the example.

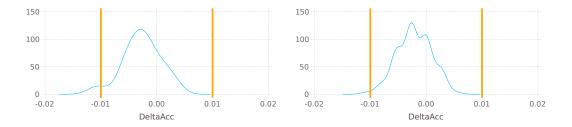
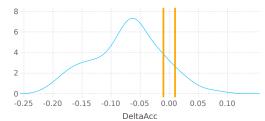


Figure 1: Density plot for the differences of accuracy between nbc and aode for the dataset hepatitis considering only 15 of the 100 data (left) or all the data (right). Left: the null hypothesis cannot be rejected (p=0.077>0.05) using half the data. Right: the null hypothesis is rejected when all the data are considered (p=0.048<0.05), despite the very small effect size.

NHST ignores magnitude and uncertainty. A very important problem with NHST is that the result of the test does not provide information about the magnitude of the effect or the uncertainty of its estimate, which are the key information we should aim at knowing.

A consequence of this limit is that: (i) a null hypothesis can be rejected despite a very small effect; (ii) a null hypothesis can be rejected even though there is a large uncertainty in the effect's magnitude and the region of uncertainty includes (or is extremely close to) zero, that is, no effect. Figure 1 (right) shows a case for which p=0.048<0.05 and the result is therefore declared to be statistically significant. However, from the density plot, it is clear that the magnitude of the effect is very small (all inside the orange vertical lines bounding the less than 1% difference of accuracy region). Thus rejecting a null hypothesis does not provide any information about the magnitude of the effect and whether or not the effect is trivial. Figure 2 shows two cases for which  $p \approx 0.001 < 0.05$  and the result is therefore declared to be statistically significant. Such p-values are similarly low, but the two cases are extremely different. For the dataset ecoli (left), the differences of accuracy are spread from 0.1 to -0.25 (the magnitude of the uncertainty is very large, about 35%), while in the second case the data are spread from 0 to -0.07, a much smaller uncertainty. Thus, rejecting a null hypothesis does not provide us with any information about the uncertainty of the estimate.



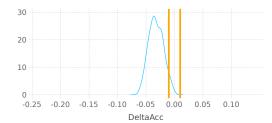


Figure 2: Density plot for the differences of accuracy (DeltaAcc) between nbc and aode for the datasets ecoli (left) and iris (right). The null hypothesis is rejected (p < 0.05) with similar p-values, even though the two cases have very different uncertainty. For ecoli, the uncertainty is very large and includes zero.

NHST yields no information about the null hypothesis. What can we say when NHST does not reject the null hypothesis?

The scientific literature contains examples of non-significant tests interpreted as evidence of no difference between the two methods/groups being compared. This is wrong since NHST cannot provide evidence in favour of the null hypothesis (see also Lecoutre and Poitevineau (2014, Sec 4.1.2.2) for further examples on this point). When NHST does not reject the null hypothesis, no conclusion can be made.

Researchers may be interested in the probability that two classifiers are equivalent. For instance, a recent paper contains the following conclusion: "there is no significant difference between (...) under the significance level of 0.05. This is quite a remarkable conclusion." This is not correct, we cannot conclude anything in this case! Consider for example Figure 3 (left), which shows the density plot of the differences of accuracy between nbc and aode for the dataset audiology. The correlated t-test gives a p-value p = 0.622 and so it fails to reject the null hypothesis. However, NHST does not allow us to reach any conclusion about the possibility that the two classifiers may actually be equivalent in this dataset, although the majority of data (density plot) seems to support this hypothesis. We tend to interpret these cases as acceptance of the null hypothesis or to even (mis)understand the p value as a "62.2% probability that the performance of the classifiers is the same". This is also evident from Figure 3 (right), where a very similar p-value, p = 0.598, corresponds to a different density plot—with more uncertainty and so less evidence in favour of the null.

There is no principled way to decide the  $\alpha$  level. In the above examples, we rejected the null hypothesis at  $\alpha=0.05$ . We understand this difference as significant, but not as significant as if we could reject it at  $\alpha=0.01$ . What is the actual meaning of this? Can we reject it at  $\alpha=0.03$ ? The  $\alpha$  level represents the crucial threshold to declare the experiment successful. With its importance, it needs to be set with care. However, since it is used to compare the meaningless p-values,  $\alpha$  is equally meaningless. By using the NHST, the researcher is forced to select an important threshold with no practical meaning. Using the customary thresholds of 0.05 and 0.01 merely allows her/him to shift the responsibility to the unsubstantiated traditional habit.

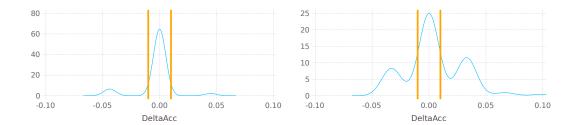


Figure 3: Density plot for the differences of accuracy between nbc and aode for the datasets audiology (left) and breast-cancer (right). In audiology, we can only say that the null hypothesis cannot be rejected p=0.622 (> 0.05). NHST does not allow the conclusion that the null hypothesis is true although practically all the data lies in the interval [-0.01, 0.01] (the differences of accuracy are less than 1%). A similar p-value corresponds to a very different situation for breast-cancer.

We pretend that  $\alpha$  is the proportion of cases in which  $H_0$  would be falsely rejected if the same experiment was repeated. This would be true if the p-value represented the probability of the null hypothesis. In reality,  $\alpha$  is the proportion of cases in which experiments yield the data that is more extreme than expected under  $H_0$ ; the actual probability of falsely rejecting  $H_0$  is also related to the probability of  $H_0$  and the data.

For  $\alpha$  to be meaningful, it would need to set the required effect size or at least the probability of the hypothesis, not the likelihood of the data.

The inference depends on the sampling intention. Consider analysing a data set of n observations with a NHST test. The sampling distribution used to determine the critical value of the test assumes that our intention was to collect exactly n observations. If the intention was different—for instance in machine learning you typically compare two algorithms on all the datasets that are available—, the sampling distribution changes to reflect the actual sampling intentions (Kruschke, 2010). This is never done, given the difficulty of formalizing one's intention and of devising an appropriate sampling distribution. This problem is thus important but generally ignored. Thus for the data set the hypothesis test (and thus the p-value) should be computed differently, depending on the intention of the person who collected the data (Kruschke, 2010).

# 3. Bayesian analysis of experimental results

There are two main Bayesian approaches for the analysis of experimental results. The first Bayesian approach, as NHST, is also based on a null value. The analyst has to set up two competing models of what values are possible. One model assumes that only the null value is possible. The alternative model assumes a broad range of other values is also possible. Bayesian inference is used to compute which model is more credible, given the data. This method is called *Bayesian model comparison* and uses so-called "Bayes factors" (Berger, 1985; Aitkin, 1991; Kass and Raftery, 1995; Berger and Pericchi, 1996).

The second Bayesian approach does not set any null value. The analyst simply has to set up a range of candidate values (prior model), including the zero effect, and use Bayesian inference to compute the relative credibilities of all the candidate values (the posterior distribution). This method is called *Bayesian estimation* (Gelman et al., 2014; Kruschke, 2015).

The choice of the method depends on the specific question that the analyst aims at answering, but in machine learning the estimation approach is usually preferable because it provides richer information to the analyst. For this reason, we will focus on the Bayesian estimation approach that hereafter we will simply call *Bayesian analysis*.

The first step in Bayesian analysis is establishing a descriptive mathematical model of the data. In a parametric model, this mathematical model is the the likelihood function that provides the probability of the observed data for each candidate value of the parameter(s)  $p(Data|\theta)$ . The second step is to establish the credibility for each value of the parameter(s) before observing data, the prior distribution  $p(\theta)$ . The third step is to use Bayes' rule to combine likelihood and prior to obtain the posterior distribution of the parameter(s) given the data  $p(\theta|Data)$ . The questions we pose in statistical analysis can be answered by querying this posterior distribution in different ways.

As a concrete example of Bayesian analysis we will compare the accuracies of two competing classifiers via cross-validation on multiple data sets (Table 1). For this purpose, we will adopt the *correlated Bayesian t-test* proposed by Corani and Benavoli (2015).

# Bayesian correlated t-test

The Bayesian correlated t-test is used for the analysis of cross-validation results on a single dataset and it accounts for the correlation due to the overlapping training sets. The test is based on the following (generative) model of the data:

$$\mathbf{x}_{n\times 1} = \mathbf{1}_{n\times 1}\mu + \mathbf{v}_{n\times 1},\tag{4}$$

where  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  is the vector of differences of accuracy,  $\mathbf{1}_{n \times 1}$  is a vector of ones,  $\mu$  is the parameter of interest (the mean difference of accuracy) and  $\mathbf{v} \sim \text{MVN}(0, \mathbf{\Sigma}_{n \times n})$  is a multivariate Normal noise with zero mean and covariance matrix  $\mathbf{\Sigma}_{n \times n}$ . The covariance matrix  $\mathbf{\Sigma}$  is characterized as follows:  $\mathbf{\Sigma}_{ii} = \sigma^2$  and  $\mathbf{\Sigma}_{ij} = \sigma^2 \rho$  for all  $i \neq j \in 1, \dots, n$ , where  $\rho$  is the correlation and  $\sigma^2$  is the variance and, therefore, the covariance matrix takes into account the correlation due to cross-validation. Hence, the likelihood model of data is

$$p(\mathbf{x}|\mu, \mathbf{\Sigma}) = \frac{\exp(-\frac{1}{2}(\mathbf{x} - \mathbf{1}\mu)^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \mathbf{1}\mu))}{(2\pi)^{n/2} \sqrt{|\mathbf{\Sigma}|}}.$$
 (5)

The likelihood (5) does not allow to estimate  $\rho$  from data, since the maximum likelihood estimate of  $\rho$  is  $\hat{\rho} = 0$  regardless the observations (Corani and Benavoli, 2015). This confirms that  $\rho$  is not identifiable: thus the Bayesian correlated t-test adopts the same heuristic  $\rho = \frac{n_{te}}{n_{tot}}$  suggested by Nadeau and Bengio (2003).

In Bayesian estimation, we aim at estimating the unknown parameters  $\mu, \nu = 1/\sigma^2$  and in particular  $\mu$ , which is the parameter of interest in the Bayesian correlated t-test. To this

end, we consider the following prior:

$$p(\mu, \nu | \mu_0, k_0, a, b) = N\left(\mu; \mu_0, \frac{k_0}{\nu}\right) G(\nu; a, b) = NG(\mu, \nu; \mu_0, k_0, a, b),$$

which is a Normal-Gamma distribution (Bernardo and Smith, 2009, Chap. 5) with parameters  $(\mu_0, k_0, a, b)$ . The Normal-Gamma prior is conjugate to the likelihood (5). If we choose the prior parameters  $\{\mu_0 = 0, k_0 \to \infty, a = -1/2, b = 0\}$  (matching prior), the resulting posterior distribution of  $\mu$  is the following Student distribution:

$$p(\mu|\mathbf{x}, \mu_0, k_0, a, b) = St\left(\mu; n - 1, \bar{\mathbf{x}}, \left(\frac{1}{n} + \frac{\rho}{1 - \rho}\right)\hat{\sigma}^2\right),\tag{6}$$

where  $\bar{x} = \frac{\sum_{i=1}^{n} x_i}{n}$  and  $\hat{\sigma}^2 = \frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n-1}$ . For these values of the prior parameters, the posterior distribution of  $\mu$  (6) coincides with the Student distribution used in the frequentist correlated t-test in (2). For instance, consider the first data set in Table 1, we have that  $\bar{x} = -0.0194$ ,  $\hat{\sigma} = 0.01583$ ,  $\rho = 1/10$ , n = 100 and so  $p(\mu | \boldsymbol{x}, \mu_0, k_0, a, b) = St(\mu; 99, -0.0194, 0.000030)$ . The output of the Bayesian analysis is the posterior of  $\mu$ ,  $p(\mu | \boldsymbol{x}, \mu_0, k_0, a, b)$ , which we can plot and query.

In (6), we have reported the posterior distribution obtained under the matching prior—for which the probability of the Bayesian correlated t-test and the p-value of the frequentist correlated t-test are numerically equivalent. Our aim is to show that although they are numerically equivalent the inferences drawn by the two approaches are very different. In particular we will show that a different interpretation of the same numerical value can completely change our prospective and allow us to make informative decisions. In other words, in this case the cassock does make the priest!

# 3.1 Comparing nbc and aode through Bayesian analysis: a colour thinking

Consider the dataset squash-unsorted, with the posterior computed by the Bayesian correlated t-test for the difference between nbc and aode, as shown in Figure 4. The vertical orange lines mark again the region corresponding to a difference of accuracy of less than 1% (we will clarify the meaning of this region later in the section). In Bayesian analysis, the experiment is summarized by the posterior distribution (in this case a Student distribution). The posterior describes the distribution of the mean difference of accuracies between the two classifiers.

By querying the posterior distribution, we can evaluate the probability of the hypothesis. We can for instance infer the probability that nbc is better/worse than aode. Formally P(nbc > aode) = 0.165 is the integral of the posterior distribution from zero to infinity or, equivalently, the posterior probability that the mean of the differences of accuracy between nbc and aode is greater than zero. P(aode > nbc) = 1 - P(nbc > aode) = 0.835 is the integral of the posterior between minus infinity and zero or, equivalently, the posterior probability that the mean of the differences of accuracy is less than zero.

Can we say anything about the probability that nbc is practically equivalent to aode? Bayesian analysis can answer this question. First, we need to define the meaning of "prac-

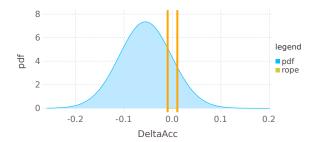


Figure 4: Posterior of the Bayesian correlated t-test for the difference between *nbc* and *aode* in the dataset *squash-unsorted*.

tically equivalent". In classification, it is sensible to define that two classifiers whose mean difference of accuracies is less that 1% are practically equivalent. The interval [-0.01, 0.01] thus defines a region of practical equivalence (rope) (Kruschke and Liddell, 2015) for classifiers.<sup>3</sup> Once we have defined a rope, from the posterior we can compute the probabilities:

- $P(nbc \ll aode)$ : the posterior probability of the mean difference of accuracies being practically negative, namely the integral of the posterior on the interval  $(-\infty, -0.01)$ .
- P(nbc = aode): the posterior probability of the two classifiers being practically equivalent, namely the integral of the posterior over the rope interval.
- $P(nbc \gg aode)$ : the posterior probability of the mean difference of accuracies being practically positive, namely the integral of the posterior on the interval  $(0.01, \infty)$ .

P(nbc = aode) = 0.086 is the integral of the posterior distribution between the vertical lines (the rope) shown in Figure 4 and it represents the probability that the two classifiers are practically equivalent. Similarly, we can compute the probabilities that the two classifiers are practically different, which are  $P(nbc \ll aode) = 0.788$  and  $P(nbc \gg aode) = 0.126$ .

The posterior also shows the uncertainty in the estimate, because the distribution shows the relative credibility of values across the continuum. One way to summarize the uncertainty is by marking the span of values that are the most credible and cover q% of the distribution (e.g., q = 90%). These are called the *High Density Intervals* (HDIs) and they are shown in Figure 5 (center) for q = 50, 60, 70, 80, 90, 95, 99.

Thus the posterior distribution equipped with rope:

- 1. estimates the posterior probability of a sensible null hypothesis (the area within the rope);
- 2. claims significant differences that also have a practical meaning (the area outside the rope);

<sup>3.</sup> In classification 1% seems to be a reasonable choice. However, in other domains a different value could be more suitable.

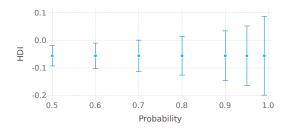


Figure 5: Posterior HDIs of the Bayesian correlated t-test for the difference between *nbc* and *aode* in the dataset *squash-unsorted*.

- 3. represents magnitude (effect size) and uncertainty (HDIs);
- 4. does not depend on the sampling intentions.

To see that, we apply Bayesian analysis to the critical cases for NHST presented in the previous section. Figure 6 shows the posterior for *hepatitis* (top), *ecoli* (left) and *iris* (right). For *hepatitis*, all the probability mass is inside the *rope* and so we can conclude that *nbc* and *aode* are practically equivalent: P(nbc = aode) = 1. For *ecoli* and *iris*, the probability mass is all in  $(-\infty, 0]$  and so nbc < aode. However, the posterior gives us more information: there is much more uncertainty in the *iris* dataset. The posteriors thus provide us the same information as the density plots shown in Figures 1–3.

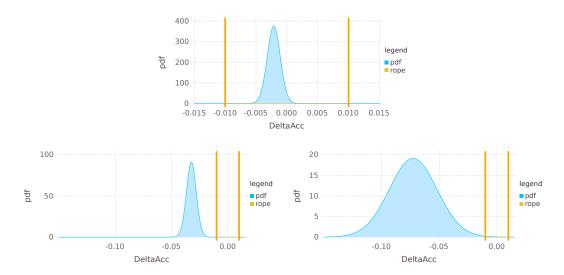


Figure 6: Posterior for nbc vs. aode in hepatitis (top), ecoli (left) and iris (right).

Consider now Figure 7; those are two examples for which we cannot clearly say whether *nbc* and *aode* are practically equivalent or practically different. We cannot decide in an

obvious way and this is evident from the posterior. Compare these figures with Figure 4, which is similar.

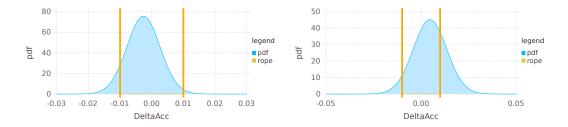


Figure 7: Posterior for nbc vs. aode in audiology (left) and breast-cancer (right).

Let us repeat the previous analysis for all 54 datasets: Figure 8 reports the posteriors for nbc versus aode in all those cases. Looking at the posteriors, we see that there are 12 cases where aode is practically better than nbc, since all the posterior is outside (to the left of) the rope (in the datasets ecoli, iris, monks1, monks, mushroom, nursey, optdigits, pasture, segment, spambase, credit, owel). There are 6 datasets where nbc and aode are practically equivalent (hayes-roth, hungarian14, hepatitis, labor, wine, yeast), since the entire posterior is inside the rope. We see also that there are no cases where nbc is practically better than aode (posterior to the right of the rope). The posteriors give us information about the magnitude of effect size, practical difference and equivalence, as well as the related uncertainty.

#### 3.2 Sensible automatic decisions

In machine learning, we often need to perform many analyses. So it may be convenient to devise a tool for automatic decision from the posterior. This means that we have to summarize the posterior in some way, with a few numbers. However, we must be aware that every time we do that we go back to that sort of black and white analysis whose limitations have been discussed before. In fact, by summarizing the posterior, we lose information, but we can do so in a conscious way. We have already explained the advantages of introducing a rope, so we can make automatic decisions based on the three probabilities  $P(nbc \ll aode)$ , P(nbc = aode) and  $P(nbc \gg aode)$ . In this way, we lose information but we introduce shades in the black and white thinking.  $P(nbc \ll aode)$ , P(nbc = aode) and  $P(nbc \gg aode)$ are probabilities and their interpretation is clear. P(nbc = aode) is the area of the posterior within the rope and represents the probability that nbc and aode are practically equivalent.  $P(nbc \ll aode)$  is the area of the posterior to the left of the rope and corresponds to the probability that nbc is practically better than aode. Finally,  $P(nbc \gg aode)$  is the area to the right of the rope and represents the probability that aode is practically better than nbc. Since these are the actual probabilities of the decisions we are interested in, in classification, we need not think in terms of Type-I errors to make decisions. We can simply make decisions using these probabilities, which have a direct interpretation—contrarily to p-values. For instance we can decide

1.  $nbc \ll aode$  if  $P(nbc \ll aode) > 0.95$ ;

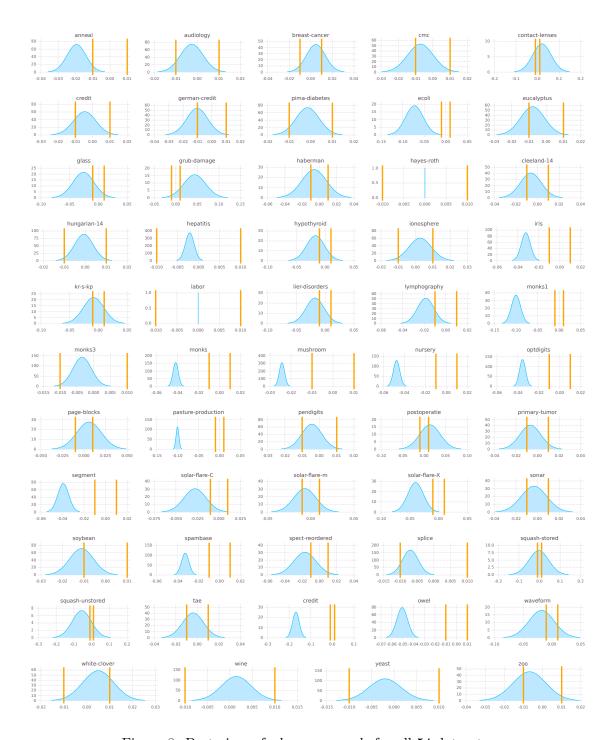


Figure 8: Posteriors of nbc versus aode for all 54 datasets.

- 2.  $nbc \gg aode$  if  $P(nbc \gg aode) > 0.95$ ;
- 3. nbc = aode if P(nbc = aode) > 0.95.

We can also decide with a probability of, for instance, 0.90 or 0.80 (if this is appropriate in the given context). Table 4 compares the Bayesian decisions based on the above rule with the NHST decisions. The NHST fails to reject the null in 35/54 datasets; Bayesian analysis declares the two classifiers equivalent in 6 of such datasets. Conversely, when NHST rejects the null (in 19/54 datasets), the Bayesian analysis declares nbc $\ll$ aode or nbc $\gg$ aode in 14 cases, nbc=aode in 1 case and no decision in 4 cases. Overall, Bayesian analysis allows us to make a decision in 6+1+14=21 datasets, while NHST makes a decision only in 19 cases.

	When NHS'	T does not re	eject the null
pair	Data sets	Bayesiar	n decision
	(out of 54)	nbc = aode	No decision
nbc-aode	35	6	29

	wnen mnsi	rejecis	ine nuii	
te				

pair	Data sets (out of 54)	Bayesian decision			
pair		nbc=aode	nbc≪aode or nbc≫aode	No decision	
nbc-aode	19	1	14	4	

Table 4: Results referring to comparisons in which the NHST rejects the null.

Sometimes also in Bayesian analysis it may be convenient to think in terms of errors. We can easily do that by defining a loss function. A loss function defines the loss we incur in making the wrong decision. The decisions we can take are  $nbc \ll aode$  (denoted as  $a_l$ ),  $nbc \gg aode$  ( $a_r$ ), nbc = aode ( $a_c$ ) or none of them ( $a_n$ ). Consider for instance the following loss matrix.

$$\begin{array}{cccc}
a_l & a_c & a_r \\
a_c & 0 & 20 & 20 \\
a_c & 20 & 0 & 20 \\
a_r & 20 & 20 & 0 \\
a_n & 1 & 1 & 1
\end{array} \tag{7}$$

The first row gives us the loss we incur in deciding  $a_l$  when  $a_l$  is true (zero loss),  $a_c$  is true (loss is 20) and  $a_r$  is true (loss is 20). Similarly for the second and third rows. The last row is the loss we incur in making no decision. The expected loss can be simply obtained by multiplying the above matrix L for the vector of posterior probabilities of nbc $\ll$ aode, nbc=aode and nbc $\gg$ aode ( $p = [p_l, p_c, p_r]^T$ ), i.e., Lp. The row of Lp corresponding to the lowest loss determines the final decision. Since  $0.05 \cdot 20 = 1$ , this leads to the same decision rule discussed previously ( $P(\cdot) > 0.95$ ).

#### 3.3 Comparing NHST and Bayesian analysis for other classifiers

In this section we extend the previous analysis to other classifiers besides nbc and aode: hidden naive Bayes (hnb), j48 decision tree (j48) and j48 grafted (j48-gr).

The aim of this section is to show that the pattern described above is general and it also holds for other classifiers. The results are presented in two tables. First we report on the cases in which NHST does not reject the null (Tab. 5). Then we report on the comparisons in which the NHST rejects the null (Tab. 6).

The NHST test does not reject the null hypothesis (Tab. 5) in 341/540 comparisons. In these cases the NHST does not make any conclusion: it cannot tell whether the null hypothesis is true<sup>4</sup> or whether is false but the evidence is too weak to reject it.

By applying the Bayesian correlated t-test with rope and taking decisions as discussed in the previous section, we can draw more informative conclusions. In 74/341=22% of the rejections failed by NHST, the posterior probability of the rope is larger than 0.95, allowing to declare that the two analyzed classifiers are practically equivalent. In the remaining cases, no conclusion can be drawn with probability 0.95.

The rope thus provides a sensible null hypothesis which can be accepted on the basis of the data. When this happens we conclude that the two classifiers are practically equivalent. This is impossible with the NHST.

When NHST does not reject the null

		12 1 acco 1000 10j		
pair	Data sets	Bayesian decision		
	(out of 54)	P(rope) > .95	No decision	
nbc-aode	35	6	29	
${ m nbc-hnb}$	30	0	30	
nbc-j48	27	2	25	
${ m nbc}{ m -j48gr}$	27	2	25	
aode-hnb	40	6	34	
aode–j48	33	6	27	
aode-j48gr	35	6	29	
hnb-j48	32	3	29	
hnb-j48gr	32	3	29	
j48–j48gr	50	40	10	
total	341	74	267	
rates		74/341 = 0.22	267/341 = 0.78	

Table 5: Results referring to comparisons in which the NHST correlated t-test does not reject the null.

Let us consider now the comparisons in which NHST claims significance (Tab. 6). There are 199 such cases. The Bayesian estimation procedure confirms the significance of 142/199 (71%) of them: in these cases it estimates either P(left) > 0.95 or P(right) > 0.95. In these cases the accuracy of the two compared classifiers are practically different. In 51/199 cases (26%) the Bayesian test does not make any conclusion. This means that a sizeable amount of probability lies within the rope, despite the statistical significance claimed by the

<sup>4.</sup> A point null is however always false, as already discussed.

<sup>5.</sup> In the comparison nbc-aode, left means  $nbc \ll aode$  and right  $nbc \gg aode$ .

When NHST rejects the null

pair	$\frac{\text{Data sets}}{\text{(out of 54)}}$	Bayesian decision			
pair		rope	difference	no decision	
nbc-aode	19	1	14	4	
$\operatorname{nbc-hnb}$	24	0	19	5	
nbc-j48	27	0	20	7	
nbc-j48gr	27	0	21	6	
aode-hnb	14	1	6	7	
aode-j48	21	1	14	6	
aode-j48gr	19	1	13	5	
hnb-j48	22	0	17	5	
hnb-j48gr	22	0	17	5	
j48–j48gr	4	2	1	1	
total	199	6	142	51	
rates		6/199 = 0.03	142/199 = 0.71	51/199=0.26	

Table 6: Results referring to comparisons in which the NHST rejects the null.

frequentist test. In the remaining cases (6/199=3%) the Bayesian test concludes that the two classifiers are practically equivalent (P(rope) > 0.95) despite the significance claimed by the NHST. In this case it draws the opposite conclusion from the NHST.

Summing up, the Bayesian test with rope is more conservative, reducing the claimed significances by 30% as compared with the NHST. The Bayesian test is thus more conservative due to the rope, which constitutes a sensible null hypothesis, while the null hypothesis of the NHST is surely wrong. However counting the detection of practically equivalent classifiers as decisions, the Bayesian test with rope takes more decisions than the NHST (222 vs. 199).

# 4. Comparing two classifiers on multiple data sets

So far we have discussed how to compare two classifiers on the same data set. In machine learning, another important problem is how to compare two classifiers on a collection of q different data sets, after having performed cross-validation on each data set.

#### 4.1 The frequentist approach

There is no direct NHST able to perform such statistical comparison, i.e., one that takes as inputs the m runs of the k-fold cross-validation differences of accuracy for each dataset and returns as output a statistical decision about which classifier is better in all the datasets. The usual NHST procedure that is employed for performing such an analysis has two steps:

1. compute the mean difference of accuracy for each dataset (averaging the differences of accuracies obtained in the m runs of the k-fold cross-validation);

2. perform a NHST to establish if the two classifiers have different performance or not based on these mean differences of accuracy.

For our case study, nbc vs. aode, the mean differences of accuracy in each dataset computed from Table 1 are shown in Table 7. We denote these measures generically with  $z = \{z_1, \ldots, z_q\}$  (in our case q = 54). The recommended NHST for this task is the signed-rank test (Demšar, 2006).

Dataset	Mean Dif.	Dataset	Mean Dif.	Dataset	Mean Dif.
anneal	-1.939	audiology	-0.261	breast-cancer	0.467
cmc	-0.719	contact-lenses	2.000	$\operatorname{credit}$	-0.464
german-credit	-1.014	pima-diabetes	-0.151	ecoli	-7.269
eucalyptus	-0.790	glass	-2.600	grub-damage	4.362
haberman	-0.614	hayes-roth	0.000	cleeland-14	-0.625
hungarian-14	-0.069	hepatitis	-0.212	hypothyroid	-1.683
ionosphere	0.267	iris	-3.242	kr-s-kp	-0.833
labor	0.000	lier-disorders	-1.762	lymphography	-1.863
monks1	-10.002	monks3	-0.343	$\operatorname{monks}$	-4.190
mushroom	-2.434	nursery	-4.747	optdigits	-3.548
page-blocks	0.583	pasture	-10.043	pendigits	-0.443
postoperatie	1.333	primary-tumor	-0.674	segment	-3.922
solar-flare-C	-2.776	solar-flare-m	-0.688	solar-flare-X	-3.996
sonar	-0.338	soybean	-1.112	spambase	-3.284
spect-reordered	-1.684	splice	-0.699	squash-stored	-0.367
squash-unstored	-5.600	tae	-0.400	$\operatorname{credit}$	-16.909
owel	-5.040	waveform	-1.809	white-clover	0.500
wine	0.143	yeast	-0.202	ZOO	-0.682

Table 7: Mean difference of accuracy (0–100) for each dataset for nbc minus aode

# Frequentist signed-rank test

The Wilcoxon signed-rank test is a non-parametric statistical hypothesis test used when comparing two paired samples. The signed-rank test assumes the observations  $z_1, \ldots, z_q$  to be i.i.d. and generated from a symmetric distribution. The test is miscalibrated if the distribution is not symmetric. A strict usage of the test should thus include first a test for symmetry. One should run the signed-rank only if "the symmetry test fails to reject the null" (note that, as we discussed before, this does not actually prove that the distribution is symmetric!). However this would make the whole testing procedure cumbersome, requiring also corrections for the test of multiple hypotheses. In the common practice thus the test for symmetry is not performed, and we follow this practice in this paper.

The test is typically used as follows. The null hypothesis is that the median of the distribution from which the  $z_i$ 's are sampled is 0; when the test rejects the null hypothesis, it claims that it is significantly different from 0. The test ranks the  $z_i$ 's according to their absolute value and then compares the ranks of the positive differences and negative

differences. The test statistic is:

$$t = \sum_{\{i: z_i \ge 0\}} r_i(|z_i|) = \sum_{1 \le i \le j \le q} t_{ij}^+,$$

where  $r_i(|z_i|)$  is the rank of  $|z_i|$  and

$$t_{ij}^{+} = \begin{cases} 1 & if \ z_i \ge -z_j, \\ 0 & otherwise. \end{cases}$$

For instance, let us consider the following two cases  $z = \{-2, -1, 4, 5\}$  or  $z = \{-1, 4, 5\}$ , then the statistic is t = 7 and, respectively, t = 5. For a large enough number of samples (e.g., q > 10), the statistic under the null hypothesis is approximately normally distributed and in this case the two-sided test is performed as follows:

$$w = \frac{t - \frac{q(q+1)}{4}}{\sqrt{\frac{q(q+1)(2q+1) - \text{tie}}{24}}},$$

$$p = 2(1 - \Phi(|w|)),$$
(8)

where p denotes the p-value computed w.r.t.  $\Phi$ , which is the cumulative distribution function of the standard Normal distribution; tie is an adjustment for ties in the data |z|, i.e.,  $z_i = -z_j$  for some i, j, required by the nonparametric test (Sidak et al., 1999; Hollander et al., 2013), while it is zero in case of no ties.

Being non-parametric, the signed-rank is robust to outliers. It assumes commensurability of differences, but only qualitatively: greater differences count more as they top the rank; yet their absolute magnitudes are ignored (Demšar, 2006).

# 4.1.1 Experimental results

If we apply this method to compare nbc vs. aode, we obtain p-value= $10^{-6}$  (the rank t=162 with no ties and w is -4.8). Since the p-value is less than 0.05, the NHST concludes that the null hypothesis can be rejected and that nbc and aode are significantly different. Table 8 reports the p-values of all comparisons of the five classifiers. The pairs nbc-aode, nbc-hnb, j48-j48gr are statistically significantly different. Again by applying this black and white mechanism, we encounter the same problems as before, i.e., we do not have any idea of the magnitude of the effect size, the uncertainty, the probability of the null hypothesis, et cetera. The density plot of the data (the mean differences of accuracy in each dataset) for nbc versus aode shows for instance that there are many datasets where the mean difference is small (close to zero), see Figure 9. Instead for j48 versus j48gr, it is clear that the difference of accuracy is very small.

#### 4.2 The Bayesian analysis approach

We will present two ways of approaching the comparison between two classifiers in multiple datasets. The first will be based on a nonparametric approach that directly extends the Wilcoxon signed-rank test. The second is a hierarchical model.

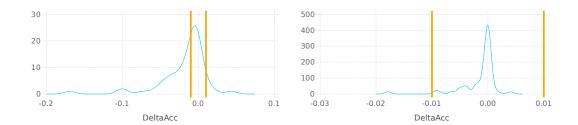


Figure 9: Density plot for nbc versus aode (left) and j48 versus j48qr (right)

Classif. 1	Classif. 2	p-value
nbc	aode	0.000
$_{ m nbc}$	$\operatorname{hnb}$	0.001
nbc	j48	0.463
$_{ m nbc}$	j48gr	0.394
aode	hnb	0.654
aode	j48	0.077
aode	j48gr	0.106
hnb	j48	0.067
hnb	j48gr	0.084
j48	j48r	0.000

Table 8: p-values for the comparison of the five classifiers.

#### 4.2.1 Nonparametric test

Benavoli et al. (2014) have proposed a Bayesian counterpart of the frequentist sign and signed-rank test, which is based on the Dirichlet process.

# Bayesian sign and signed-tank test

Let z denote the scalar variable of interest and  $z = \{z_1, \ldots, z_q\}$  denotes a vector of i.i.d. observations of z. To derive the Bayesian sign and signed-rank test, we assume a Dirichlet Process (DP) prior on the probability distribution of z. A DP is a distribution over probability distributions such that marginals on finite partitions are Dirichlet distributed. Like the Dirichlet distribution, the DP is therefore completely defined by two parameters: the prior strength s > 0 and the prior mean that, for the DP, is a probability measure  $G_0$  on z. If we choose  $G_0 = \delta_{z_0}$ , i.e., a Dirac's delta centred on the pseudo-observation  $z_0$ , the posterior probability density function of Z has this simple expression:

$$p(z) = w_0 \delta_{z_0}(z) + \sum_{j=1}^n w_j \delta_{z_j}(z), \quad (w_0, w_1, \dots, w_n) \sim Dir(s, 1, \dots, 1),$$
(9)

i.e., it is a mixture of Dirac's deltas centred on the observations  $z_j$  for j = 1, ..., q and on the prior pseudo-observation  $z_0$ , whose weights are Dirichlet distributed with parameters (s, 1, ..., 1). We can think about (9) as a hierarchical model: p(z) depends on the weights

w that are Dirichlet distributed. The model (9) is therefore a posterior distribution of the probability distribution of z and encloses all the information we need for the experimental analysis. We can summarize it in different way. If we compute:

$$\theta_l = P(z < -r) = \sum_{i=0}^{q} w_i I_{(-\infty, -r)}(z_i),$$

$$\theta_e = P(|z| \le r) = \sum_{i=0}^{q} w_i I_{[-r, r]}(z_i),$$

$$\theta_r = P(z > r) = \sum_{i=0}^{q} w_i I_{(r, \infty)}(z_i),$$

where the indicator  $I_A(z) = 1$  if  $z_0 \in A$  and zero otherwise, then we obtain a Bayesian version of the **sign test** that also accounts for the rope [-r, r]. In fact,  $\theta_l, \theta_e, \theta_r$  are respectively the probabilities that the mean difference of accuracy is in the interval  $(-\infty, -r)$ , [-r, r], or  $(r, \infty)$ . Since  $(w_0, w_1, \ldots, w_n) \sim Dir(s, 1, \ldots, 1)$ , it can easily be shown that

$$\theta_l, \theta_e, \theta_r \sim Dirichlet(n_l + sI_{(-\infty, -r]}(z_0), n_e + sI_{[-r, r]}(z_0), n_r + sI_{[r, \infty)}(z_0)),$$
 (10)

where  $n_l$  is the number of observations  $z_i$  that fall in  $(-\infty, -r]$ ,  $n_e$  is the number of observations  $z_i$  that fall in [-r, r] and  $n_r$  is the number of observations  $z_i$  that fall in  $[r, \infty)$ , obviously  $n_l + n_e + n_r = q$ . If we neglect  $sI_{(-\infty, -r]}(z_0), sI_{[-r,r]}(z_0), sI_{[r,\infty)}(z_0)$ , then (10) says that the posterior probability of  $\theta_l$ ,  $\theta_e$ ,  $\theta_r$  is Dirichlet distributed with parameters  $(n_l, n_e, n_r)$ . The terms  $sI_{(-\infty, -r]}(z_0), sI_{[-r,r]}(z_0), sI_{[r,\infty)}(z_0)$  are due to the prior. Therefore, to fully specify the Bayesian sign test, we must choose the value of the prior strength s and where to place the pseudo-observation  $z_0$  in  $(-\infty, -r]$  or [-r, r] or  $[r, \infty)$ . We will return to this choice in Section 4.3. Instead, if we compute

$$\theta_{l} = \sum_{i=0}^{q} \sum_{j=0}^{q} w_{i} w_{j} I_{(-\infty,-2r)}(z_{j} + z_{i}),$$

$$\theta_{e} = \sum_{i=0}^{q} \sum_{j=0}^{q} w_{i} w_{j} I_{[-2r,2r]}(z_{j} + z_{i}),$$

$$\theta_{r} = \sum_{i=0}^{q} \sum_{j=0}^{q} w_{i} w_{j} I_{(2r,\infty)}(z_{j} + z_{i}),$$
(11)

then we derive a Bayesian version of the **signed rank test** (Benavoli et al., 2014) that also accounts for the rope [-r, r]. This time the distribution of  $\theta_l, \theta_e, \theta_r$  has not a simple closed form but we can easily compute it by Monte Carlo sampling the weights  $(w_0, w_1, \ldots, w_n) \sim Dir(s, 1, \ldots, 1)$ . Also in this case we must choose  $s, z_0$ , see Section 4.3.

It should be observed that the Bayesian signed-rank test does not require the symmetry assumption about the distribution of the observations  $z_i$ . This test works also in case the distribution is asymmetric thanks to the Bayesian estimation approach (i.e., it estimates

the distribution from data). This is another advantage of the Bayesian estimation approach w.r.t. the frequentist null hypothesis tests (Benavoli et al., 2014).

#### 4.2.2 Experiments

Let us start by comparing nbc vs. aode by means of the Bayesian sign-rank tests without rope (r=0). Hereafter we will choose the prior parameter of the Dirichlet as s=0.5 and  $z_0=0$ ; we will return to this choice in Section 4.3. Since without rope  $\theta_r=1-\theta_l$ , we have only reported the posterior of  $\theta_l$  (denoted as "Pleft") that represents the probability that aode is better than nbc. The samples of the posteriors are shown in Figure 10: this is simply the histogram of 150′000 samples of  $\theta_l$  generated according to (11). For all samples, it results in  $\theta_l$  greater than 0.5 and so  $\theta_r=1-\theta_l$ . So we can conclude with probability  $\approx 1$  that aode is better than nbc. We can in fact think about the comparison of two classifiers as the inference on the bias  $(\theta_l)$  of a coin. In this case, all the 150′000 sampled coins from the posterior have a bias that is greater than 0.5 and, therefore, all the coins are always biased towards aode (which is then preferable to nbc).

This conclusion is in agreement with that derived by the frequentist sign-rank test (very small p-value, see Table 8).

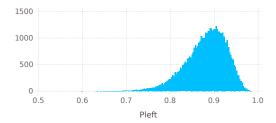


Figure 10: Posterior for nbc vs. aode for the Bayesian sign-rank test.

The introduction of the rope partially changes the previous conclusion. A way to visualize the posterior of  $\theta_l$ ,  $\theta_e$ ,  $\theta_r$  in this case, is by plotting the 150'000 Monte Carlo samples of these probabilities in barycentric coordinates: each trinomial vector of probabilities is a point in the simplex having vertices  $\{(1,0,0),(0,1,0),(0,0,1)\}$ . The three vertices are respectively denoted as "aode", "rope" and "nbc" and represent decisions with certainty in favour of "aode", "rope" and, respectively, "nbc". Figure 11 reports the simplex as well as the two-dimensional projections of the posterior for the Bayesian sign-rank test. In particular Figure 11 reports the marginal of the posterior distribution of "aode" vs. "rope" (left); the marginal of the posterior distribution "nbc" vs. "rope" (right). From these two figures we can deduce the other marginal since  $\theta_l + \theta_e + \theta_r = 1$ . Finally, Figure 11 (bottom) reports the joint of the three variables in barycentric coordinates (we are again exploiting the fact that  $\theta_l + \theta_e + \theta_r = 1$ ). In particular, Figure 11 (bottom) reports the samples from the posteriors (cloud of points), the simplex (the large orange triangle) and three regions (in orange) that are limited by the level curves:  $\theta_i \geq \max(\theta_i, \theta_k)$  with  $i \neq j \neq k$  (hypothesis i is more probable than both hypotheses i, k together). For instance, the region at the bottom-right of the triangle is relative to the case where "aode" is more probable than "rope" and "nbc" together; the region at the top of the triangle represents the case where "rope" is more

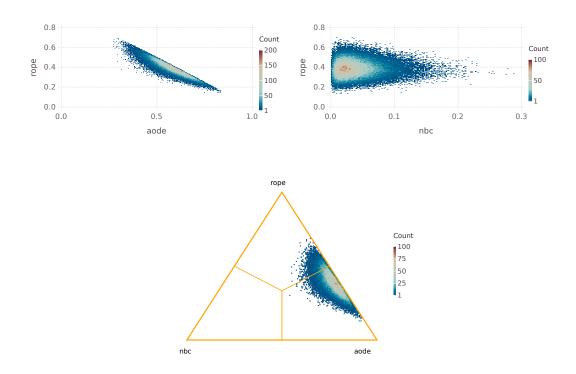


Figure 11: Posterior for nbc vs. aode for the Bayesian sign-rank test.

probable than "aode" and "nbc" together; the region at the left of the triangle corresponds to the case where "nbc" is more probable than "aode" and "rope" together. Hence, if all the points fall inside one of these three regions, we conclude that such hypothesis is true with probability  $\approx 1$ . Looking at Figure 11 (bottom), it is evident that the majority of cases support aode more than rope and definitively more than nbc. We can quantify this numerically by counting the number of points that fall in the three regions, see first row in Table 9. aode is better in 90% of cases, while rope is selected in the remaining 10%. We can therefore conclude with probability 90% that aode is practically better than nbc. Table 9 reports also these probabilities for the other comparisons of classifiers computed using 150'000 Monte Carlo samples. We conclude that hnb is practically better than nbc with probability 0.999; aode and hnb are equivalent with probability 0.95; aode is better than j48 and j48gr with probability 0.9; hnb is better than j48 and j48gr with probability greater than 0.95 and finally j48 and j48gr are practically equivalent. These conclusions are in agreement with the data.

The computational complexity of the Bayesian sign-rank test is low. The comparison of two classifiers (based on 150'000 samples) takes less than one second on a standard computer.

# 4.3 Choice of the prior

In the previous section, we have selected the prior parameters of the Dirichlet process as s = 0.5 and  $z_0 = 0$ . In terms of rank, this basically means that the prior strength is

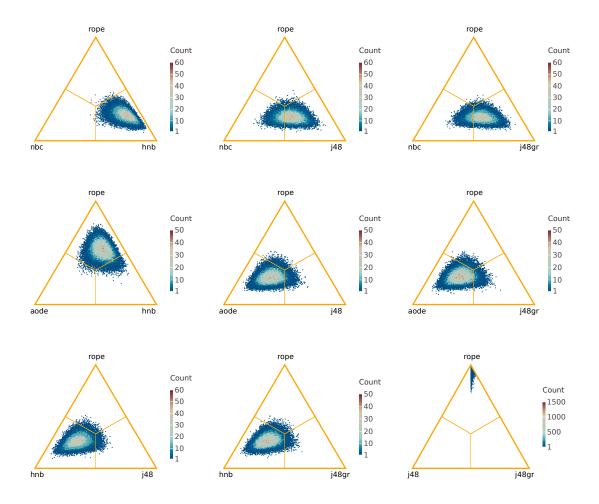


Figure 12: Posterior for *nbc* vs. *aode* for Bayesian sign-rank test.

equivalent to that of one pseudo-observation that is located inside the rope (Benavoli et al., 2014). How are the inferences sensitive to this choice? For instance, we can see how the probabilities on Table 9 would change based on  $z_0$ . We have considered two extreme cases  $z_0 = -\infty$  and  $z_0 = \infty$  and reported these probabilities in Table 10 and 11 (this is an example of robust Bayesian analysis (Berger et al., 1994)). It is evident that the position of  $z_0$  has only a minor effect on the probabilities. We could have performed this analysis jointly by considering all the possible Dirichlet process priors obtained by varying  $z_0 \in \mathbb{R}$ . This set of Dirichlet priors is called "Imprecise Dirichlet Process" (IDP). IDP allows us to start the inference with very weak prior assumptions, much in the direction of letting data speak for themselves. More details about the properties of IDP and the choice of the prior can be found in Benavoli et al. (2014, 2015b); Walley (1996).

Classif. 1	Classif. 2	left	rope	$\operatorname{right}$
nbc	aode	0.000	0.103	0.897
nbc	hnb	0.000	0.001	0.999
nbc	j48	0.228	0.004	0.768
nbc	j48gr	0.182	0.002	0.815
aode	hnb	0.001	0.956	0.042
aode	j48	0.911	0.026	0.063
aode	j48gr	0.892	0.035	0.073
hnb	j48	0.966	0.015	0.019
hnb	j48gr	0.955	0.020	0.025
j48	j48gr	0.000	1.000	0.000

Table 9: Probabilities for the ten comparisons of classifiers. Left and right refer to the columns Classif. 1 (left) and Classif. 2 (right).

Classif. 1	Classif. 2	left	rope	$\operatorname{right}$
nbc	aode	0.000	0.112	0.888
nbc	hnb	0.000	0.001	0.999
nbc	j48	0.262	0.004	0.734
nbc	j48gr	0.213	0.003	0.784
aode	$\operatorname{hnb}$	0.002	0.961	0.037
aode	j48	0.922	0.024	0.053
aode	j48gr	0.906	0.033	0.061
hnb	j48	0.971	0.014	0.016
hnb	j48gr	0.961	0.018	0.021
j48	j48gr	0.000	1.000	0.000

Table 10: Probabilities for the ten comparisons of classifiers with  $z_0 = \infty$ . Left and right refer to the columns Classif. 1 (left) and Classif. 2 (right).

# 4.3.1 Hierarchical models

In Section 3 we have presented the Bayesian correlated t-test that is used for the analysis of cross-validation results on a single dataset. In particular, it makes inference about the mean difference of accuracy between two classifiers in the *i*-th dataset ( $\mu_i$ ) by exploiting three pieces of information: the sample mean ( $\bar{x}_i$ ), the variability of the data (sample standard deviation  $\hat{\sigma}_i$ ) and the correlation due to the overlapping training set ( $\rho$ ). This test can only be applied to a single dataset. We have already discussed the fact that there is no direct NHST able to extend the above statistical comparison to multiple datasets, i.e., that takes as inputs the m runs of the k-fold cross-validation results for each dataset and returns as output a statistical decision about which classifier is better in all the datasets. The usual NHST procedure that is employed for performing such analysis has two steps: (1) compute

Classif. 1	Classif. 2	left	rope	$\operatorname{right}$
nbc	aode	0.000	0.096	0.904
nbc	hnb	0.000	0.001	0.999
nbc	j48	0.201	0.004	0.795
nbc	j48gr	0.159	0.002	0.839
aode	hnb	0.001	0.950	0.049
aode	j48	0.892	0.028	0.080
aode	j48gr	0.872	0.037	0.091
hnb	j48	0.957	0.017	0.027
hnb	j48gr	0.944	0.022	0.034
j48	j48gr	0.000	1.000	0.000

Table 11: Probabilities for the ten comparisons of classifiers with  $z_0 = -\infty$ . Left and right refer to the columns Classif. 1 (left) and Classif. 2 (right).

the mean difference of accuracy for each dataset  $\bar{x}_i$ ; (2) perform a NHST to establish if the two classifiers have different performance or not based on these mean differences of accuracy. This discards two pieces of information: the correlation  $\rho$  and sample standard deviation  $\hat{\sigma}_i$  in each dataset. The standard deviation is informative about the accuracy of  $\bar{x}_i$  as an estimator of  $\mu_i$ . The standard deviation can largely vary across data sets, as a result of each data set having its own size and complexity. The aim of this section is to present an extension of the Bayesian correlated t-test that is able to make inference on multiple datasets and at the same time to account for all the available information (mean, standard deviation and correlation). In Bayesian estimation, this can be obtained by defining a hierarchical model (Corani et al., 2017). Hierarchical models are among the most powerful and flexible tools in Bayesian analysis.

#### Bayesian hierarchical correlated t-test

The hierarchical correlated t-test is based on following hierarchical probabilistic model:

$$\mathbf{x}_i \sim MVN(\mathbf{1}\mu_i, \mathbf{\Sigma_i}),$$
 (12)

$$\mu_1...\mu_q \sim t(\mu_0, \sigma_0, \nu),$$
 (13)

$$\sigma_1...\sigma_q \sim \text{unif}(0,\bar{\sigma}).$$
 (14)

Equation (12) models the fact that the cross-validation measures  $\mathbf{x}_i = \{x_{i1}, x_{i2}, \dots, x_{in}\}$  of the i-th data set are jointly multivariate-normal distributed with the same mean  $(\mu_i)$ , same variance  $(\sigma_i)$  and correlation  $(\rho)$ . In fact, it states that, for each dataset i,  $\mathbf{x}_i$  is multivariate normal with mean  $\mathbf{1}\mu_i$  (where  $\mathbf{1}$  is a vector of ones) and covariance matrix  $\mathbf{\Sigma}_i$  defined as follows: the diagonal elements are  $\sigma_i^2$  and the out-of-diagonal elements are  $\rho\sigma_i^2$ , where  $\rho = \frac{n_{te}}{n_{tr}}$ . This model is the same we discussed in Section 3. Equation (13) models the fact that the mean difference of accuracies in the single datasets,  $\mu_i$ , depends on  $\mu_0$  that is the average difference of accuracy between the two classifiers on the population of data sets. This is the quantity we aim at estimating. Equation (13) assumes the  $\mu_i$ 's to

be drawn from a high-level Student distribution with mean  $\mu_0$ , variance  $\sigma_0^2$  and degrees of freedom  $\nu$ . The choice of a Student distribution at this level of the hierarchical model enables the model to robustly deal with data sets whose  $\mu_i$ 's are far away from the others (Gelman et al., 2014; Kruschke, 2013). Moreover the heavy tails of the Student make more cautious the conclusions drawn by the model.

The hierarchical model assigns to the i-th data set its own standard deviation  $\sigma_i$ , assuming the  $\sigma_i$ 's to be drawn from a common distribution, see Equation (14). In this way it realistically represents the fact the estimates referring to different data sets data sets have different uncertainty. The high-level distribution of the  $\sigma_i$ 's is unif $(0, \bar{\sigma})$ , as recommended by Gelman (2006), as it yields inferences which are insensitive to  $\bar{\sigma}$ , if  $\bar{\sigma}$  is large enough. To this end we set  $\bar{\sigma} = 1000 \cdot \bar{s}$  (Kruschke, 2013), where  $\bar{s} = \sum_{i}^{q} \hat{\sigma}_i/q$ .

We complete the model with the prior on the parameters  $\delta_0$ ,  $\sigma_0$  and  $\nu$  of the high-level distribution. We assume  $\delta_0$  to be uniformly distributed within 1 and -1. This choice works for all the measures bounded within  $\pm 1$ , such as accuracy, AUC, precision and recall. Other type of indicators might require different bounds.

For the standard deviation  $\sigma_0$  we adopt the prior  $unif(0, \bar{s_0})$ , with  $\bar{s_0} = 1000s_{\bar{x}}$ , where  $s_{\bar{x}}$  is the standard deviation of the  $\bar{x}_i$ 's.

As for the prior  $p(\nu)$  on the degrees of freedom, there are two proposals in the literature. Kruschke (2013) proposes an exponentially shaped distribution which balances the prior probability of nearly normal distributions ( $\nu > 30$ ) and heavy tailed distributions ( $\nu < 30$ ). We re-parameterize this distribution as a Gamma( $\alpha,\beta$ ) with  $\alpha=1, \beta=0.0345$ . Juárez and Steel (2010) proposes instead  $p(\nu)=\text{Gamma}(2,0.1)$ , assigning larger prior probability to normal distributions.

We have no reason for preferring a prior over another, but the hierarchical model shows some sensitivity on the choice of  $p(\nu)$ . We model this uncertainty by representing the coefficients  $\alpha$  and  $\beta$  of the Gamma distribution as two random variables (hierarchical prior). In particular we assume  $p(\nu) = \text{Gamma}(\alpha, \beta)$ , with  $\alpha \sim \text{unif}(\underline{\alpha}, \bar{\alpha})$  and  $\beta \sim \text{unif}(\underline{\beta}, \bar{\beta})$ , setting  $\underline{\alpha}$ =0.5,  $\bar{\alpha}$ =5,  $\underline{\beta}$ =0.05,  $\bar{\beta}$ =0.15. The simulations in (Corani et al., 2017) show that the inferences of the model are stable with respect to perturbations of  $\underline{\alpha}$ ,  $\bar{\alpha}$ ,  $\underline{\beta}$ , and  $\bar{\beta}$ , and that the resulting hierarchical generally fits well the experimental data.

These considerations are reflected by the following probabilistic model:

$$\nu \sim Ga(\alpha, \beta),$$
 (15)

$$\alpha \sim \operatorname{unif}(\underline{\alpha}, \overline{\alpha}),$$
 (16)

$$\beta \sim \operatorname{unif}(\beta, \overline{\beta}),$$
 (17)

$$\mu_0 \sim \text{unif}(-1, 1),\tag{18}$$

$$\sigma_0 \sim \text{unif}(0, \bar{\sigma_0}).$$
 (19)

We want to make inference about the  $\mu_i$ 's and  $\mu_0$ . Such inferences are computed by marginalizing out the  $\sigma_i$ 's, and thus accounting for the different uncertainty which characterizes each data set. This characteristic is unique among the methods discussed so far. Computations in hierarchical models are obtained by Markov-Chain Monte Carlo sampling.

A further merit of the hierarchical model is that it jointly estimates the  $\mu_i$ 's while the existing methods estimate independently the difference of accuracy on each data set using the  $\bar{x}_i$ 's. The consequence of the joint estimation performed by the hierarchical model is that shrinkage is applied to the  $\bar{x}_i$ 's. The hierarchical model thus estimates the  $\mu_i$ 's more accurately than the  $\bar{x}_i$ 's adopted by the other tests. This result is valid under general assumptions, such as a severe misspecification between the high-level distributions of the true generative model and of the fitted model (Corani et al., 2017). By applying the rope on the posterior distribution of the  $\mu_i$ 's and the  $\mu_0$  in a similar way to what discussed for the Bayesian correlated t-test, the model is able to detect equivalent classifiers and to claim significances that have a practical impact.

#### 4.3.2 Experiments

In the experiments, we have computed the posterior of  $\mu_0$ ,  $\sigma_0$ ,  $\nu$  for the ten pairwise comparisons between the classifiers nbc, aode, hnb, j48 and j48gr. As inference we have computed the prediction on the next (unseen) dataset, which is formally equivalent to the inference computed by the Bayesian signed-rank test. For instance, for nbc vs. aode, we have computed the probabilities that in the next dataset nbc is better than aode ( $\theta_r$ ), nbc is equivalent to aode ( $\theta_e$ ), aode is better than nbc ( $\theta_l$ ). This is the procedure we have followed:

- 1. we have sampled  $\mu_0, \sigma_0, \nu$  from the posteriors of these parameters;
- 2. for each sample of  $\mu_0, \sigma_0, \nu$  we have defined the posterior of the mean difference of accuracy on the next dataset, i.e.,  $t(\mu_{next}; \mu_0, \sigma_0, \nu)$ ;
- 3. from  $t(\mu_{next}; \mu_0, \sigma_0, \nu)$  we have computed the probabilities  $\theta_l$  (integral on  $[-\infty, r]$ ),  $\theta_e$  (integral on [-r, r]) and  $\theta_r$  (integral on  $[r, \infty)$ ).

We have repeated this procedure 4'000 times, obtaining 4'000 samples of  $(\theta_l, \theta_e, \theta_r)$  and the results are shown in Figure 13. The results are quite in agreement with those of the Bayesian signed-rank test. For instance, it is evident that aode is clearly better than nbc. We can quantify this numerically by counting the number of points that fall in the three regions (see the first row in Table 12). aode is better in almost 100% of cases. Table 12 reports also these probabilities for the other comparisons of classifiers computed using 4'000 Monte Carlo samples. By comparing Tables 12 and 9, we can see that the two tests are substantially in agreement apart from differences in aode vs. j48 and j48gr. The hierarchical test is taking into account of all available infromation: the sample mean  $(\bar{x}_i)$ , the variability of the data (sample standard deviation  $\hat{\sigma}_i$ ) and the correlation due to the overlapping training set  $(\rho)$ , while the Bayesian signed rank only considers  $\bar{x}_i$ . Therefore, when the two tests differ substantially, it means that there is substantial variability of the cross-validation estimate.

We have implemented the hierarchical model in Stan (http://mc-stan.org) (Carpenter et al., 2016), a language for Bayesian inference. The analysis of the results of 10 runs of 10-fold cross-validation on 54 data sets (that means a total of 5400 observations) takes about three minutes on a standard computer.

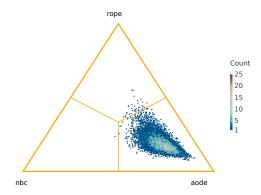


Figure 13: Posterior of *nbc* versus *aode* for all 54 datasets.

Classif. 1	Classif. 2	left	rope	$\operatorname{right}$
nbc	aode	0	0.28	0.72
nbc	hnb	0	0	1
nbc	j48	0.2	0.01	0.79
nbc	j48gr	0.15	0.01	0.84
aode	$\operatorname{hnb}$	0	1	0
aode	j48	0.46	0.51	0.03
aode	j48gr	0.41	0.56	0.03
$\operatorname{hnb}$	j48	0.91	0.07	0.02
hnb	j48gr	0.92	0.05	0.03
j48	j48gr	0	1	0

Table 12: Probabilities for the ten comparisons of classifiers. Left and right refer to the columns Classif. 1 (left) and Classif. 2 (right).

#### 4.4 Choice of the hyper-priors parameters

The choice of of the hyper-priors parameters can be critical in Bayesian hierarchical models. We have conducted a sensitivity analysis by using different constants in the top-level gamma and uniform distributions, to check whether they have any notable influence on the resulting posterior distribution. Whether all gamma/uniform distributions are assumed, the results are essentially identical. More details about this sensitivity analysis are reported in Corani et al. (2017). This means that for the hierarchical model inferences and decisions are stable w.r.t. the choice of the hyper-parameters.

# 4.5 Bayesian signed rank or hierarchical model?

So far we have presented two methods for comparing two classifiers on multiple datasets: Bayesian signed-rank and hierarchical model. Which one should we use for comparing classifiers? In our opinion, the hierarchical model is preferable because it takes as inputs the m runs of the k-fold cross-validation results for each dataset and so it makes inference about the mean difference of accuracy between two classifiers in the i-th dataset  $(\mu_i)$  by exploiting all available information: the sample mean  $(\bar{x}_i)$ , the variability of the data (sample standard deviation  $\hat{\sigma}_i$ ) and the correlation due to the overlapping training set  $(\rho)$ . Conversely, the Bayesian signed rank only considers  $\bar{x}_i$ . On the other hand, the hierarchical model is slower than the Bayesian signed rank. In machine learning, we often need to run statistical tests hundreds of times for instance for features selection or algorithms racing and, in this case, it is more convenient to use a light test as the Bayesian signed rank.

# 5. Comparisons of multiple classifiers

Another important problem with NHST is the issue of multiple hypothesis testing. Considering the results in Table 8, which reports the p-values for the comparison of the five classifiers obtained by the Wilcoxon signed-rank test. From the p-values, we concluded that "nbc was found significantly better than aode and hnb, and algorithms j48 and j48grwere significantly different, while there were no significant differences between other pairs". When many tests are made, the probability of making at least one Type 1 error in any of the comparisons increases. One of the most popular fixes to this problem is the Bonferroni correction. The Bonferroni correction adjusts the p-value at which a test is evaluated for significance according to the number of tests being performed. More specifically, the adjusted p-value is calculated as the original p-value divided by the number of tests being performed. Implicitly, Bonferroni's correction assumes that these test statistics are independent. So in our current example an overall desired significance level of 0.05 would translate into individual tests each using a p-value threshold of 0.05/10 = 0.005 (we are performing 10 comparisons). In this case, this would not change our previous sections, since all significant p-values were less than 0.005. The Bonferroni correction reduces false rejections but it also increases the number of instances in which the null is not rejected when actually it should have been. Thus, the Bonferroni adjustment can reduce the power to detect an important effect. Motivated by this issue of the Bonferroni correction, researchers have proposed alternative procedures. The goal of these methods typically is to reduce the family-wise error rate (that is, the probability of having at least one false positive) without sacrificing power too much. A natural way to achieve this is by considering the dependence across tests (Westfall et al., 1993).

We have already discussed in Section 2 the pitfalls of NHST Type I error thinking. Type I error underlies these corrections and, therefore, corrections inherit all its problems. The most critical one is that the correction factor depends on the way the analyst intends to conduct the comparison. For instance, the analyst may want to compare nbc with the other four classifiers (in this case the Bonferroni correction would be 0.05/4 = 0.0125—he/she is conducting only four comparisons) or to perform all the ten comparisons (0.05/10) and so on. This creates a problem because two analysts can c draw different conclusions from the same data because of the variety of comparisons that they made. Another important issue with the multiple-comparison procedure based on mean-ranks test is described in Benavoli et al. (2016).

How do we manage the problem of multiple hypothesis testing in Bayesian analysis? Paraphrasing Gelman et al. (2012): "in Bayesian analysis we usually do not have to worry about multiple comparisons. The reason is that we do not worry about Type I error, because the null hypothesis is hardly believable to be true." How does Bayesian Analysis mitigate false alarms? Gelman et al. (2012) suggest using multilevel analysis (in our case a hierarchical Bayesian model on multiple classifiers). Multilevel models perform partial pooling; they shift estimates toward each other. This means that the comparisons of the classifiers are more conservative, in the sense that intervals for comparisons are more likely to include zero. This may be a direction to pursue in future research. In this paper, we instead mitigate false alarms through the rope. The rope mitigates false alarms because it decreases the asymptotic false alarm rate (Kruschke, 2013).

#### 6. Software and available Bayesian tests

All the tests that we have presented in this paper are available in *R* and *Python* code at https://github.com/BayesianTestsML/tutorial/.

Moreover, the code that is necessary to replicate all the analyses we performed in this paper is also available at the above URL in form of *Ipython* notebooks (implemented in *Python* and *Julia*). The software is open source and that can be freely used, changed, and shared (in modified or unmodified form).

Machine learning researchers may be interested in using other Bayesian tests besides the ones we have discussed in this paper. General Bayesian parametric tests can be found in Kruschke (2015) (together with R code) and also in Gelman et al. (2013). We have specialized some of these tests to the case of correlated data, such as the Bayesian correlated t-test (Corani and Benavoli, 2015) discussed in Section 3. We have also implemented several Bayesian nonparametric tests for comparing algorithms: Bayesian rank test (Benavoli et al., 2015b), Friedman test (Benavoli et al., 2015a) and tests that account for censored data (Mangili et al., 2015). Finally, we have developed an extension of the Bayesian sign test to compare algorithms taking into account multiple measures at the same time (accuracy and computational time for instance) (Benavoli and Campos, 2015). For the analysis of multiple data sets, another approach has been proposed by Lacoste et al. (2012) that models each data set as an independent Bernoulli trial. The two possible outcomes of the Bernoulli trial are the first classifier being more accurate than the second or vice versa. This approach yields the posterior probability of the first classifier being more accurate than the second classifier on more than half of the q data sets. A shortcoming is that its conclusions apply only to the q available data sets without generalizing to the whole population of data sets

# 7. Conclusions

We discourage the use of frequentist null hypothesis significance tests (NHST) in machine learning and, in particular, for comparison of the performance of classifiers. In this, we follow the current trends in other scientific areas. For instance, the journal of Basic and Applied Social Psychology, has banned the use of NHSTs and related statistical procedures (Trafimow and Marks, 2015). We believe that also in machine learning is time to move on from NHST and p-values.

In this paper, we have discussed how Bayesian analysis can be employed instead of NHST. In particular, we have presented three Bayesian tests: Bayesian correlated t-test, Bayesian signed rank test and a Bayesian hierarchical model that can be used for comparing the performance of classifiers and that solve the drawbacks of the frequentist tests. All the code of these tests is freely available and so researchers can already use these tests for their analysis. In this paper, we have mainly discussed the use of Bayesian tests for comparing the performance of algorithms. However, in machine learning, NHST statistical tests are also employed inside the algorithms. For instance, nonparametric tests are used in racing algorithms, independence tests are used to learn the structure of Bayesian networks, etcetera. Bayesian tests can be used to replace all NHST tests because of their advantages (for instance, Bayesian tests can assess whether two algorithms are similar through the use of the rope (Benavoli et al., 2015a)).

# Acknowledgements

Research partially supported by the Swiss NSF grant no. IZKSZ2\_162188.

#### References

- Murray Aitkin. Posterior Bayes factors. Journal of the Royal Statistical Society. Series B (Methodological), pages 111–142, 1991.
- Alessio Benavoli and Cassio P. Campos. Advanced Methodologies for Bayesian Networks: Second International Workshop, AMBN 2015, Yokohama, Japan, November 16-18, 2015. Proceedings, chapter Statistical tests for joint analysis of performance measures. Springer International Publishing, Cham, 2015.
- Alessio Benavoli, Francesca Mangili, Giorgio Corani, Marco Zaffalon, and Fabrizio Ruggeri. A Bayesian Wilcoxon signed-rank test based on the Dirichlet process. In *Proceedings of the 30th International Conference on Machine Learning (ICML 2014)*, pages 1–9, 2014.
- Alessio Benavoli, Giorgio Corani, Francesca Mangili, and Marco Zaffalon. A Bayesian non-parametric procedure for comparing algorithms. In *Proceedings of the 31th International Conference on Machine Learning (ICML 2015)*, pages 1–9, 2015a.
- Alessio Benavoli, Francesca Mangili, Fabrizio Ruggeri, and Marco Zaffalon. Imprecise Dirichlet process with application to the hypothesis test on the probability that X≤Y. *Journal of Statistical Theory and Practice*, 9(3):658–684, 2015b.
- Alessio Benavoli, Giorgio Corani, and Francesca Mangili. Should we really use post-hoc tests based on mean-ranks? *Journal of Machine Learning Research*, 17(5):1–10, 2016.
- James O. Berger. Statistical Decision Theory and Bayesian Analysis. Springer Series in Statistics, New York, 1985.
- James O. Berger and Luis R Pericchi. The intrinsic Bayes factor for model selection and prediction. *Journal of the American Statistical Association*, 91(433):109–122, 1996.

- James O. Berger and Thomas Sellke. Testing a point null hypothesis: The irreconcilability of p-values and evidence. *Journal of the American statistical Association*, 82(397):112–122, 1987.
- James O. Berger, E. Moreno, L. R. Pericchi, M. J. Bayarri, Bernardo, et al. An overview of robust Bayesian analysis. *Test*, 3(1):5–124, 1994.
- José M Bernardo and Adrian FM Smith. *Bayesian theory*, volume 405. Wiley Chichester, 2009.
- Christopher Bishop. Pattern Recognition and Machine Learning. Springer, 2007.
- Remco R Bouckaert. Choosing between two learning algorithms based on calibrated tests. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 51–58, 2003.
- Bob Carpenter, Daniel Lee, Marcus A Brubaker, Allen Riddell, Andrew Gelman, Ben Goodrich, Jiqiang Guo, Matt Hoffman, Michael Betancourt, and Peter Li. Stan: A probabilistic programming language. *Journal of Statistical Software*, in press, 2016.
- Giorgio Corani and Alessio Benavoli. A Bayesian approach for comparing cross-validated algorithms on multiple data sets. *Machine Learning*, 100(2):285–304, 2015. doi: 10.1080/s10994-015-5486-z.
- Giorgio Corani, Alessio Benavoli, Janez Demsar, Francesca Mangili, and Marco Zaffalon. Statistical comparison of classifiers through Bayesian hierarchical modelling. *Machine Learning in press*, 2017. doi: 10.1007/s10994-017-5641-9.
- Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, 2006.
- Janez Demšar. On the appropriateness of statistical tests in machine learning. In Workshop on Evaluation Methods for Machine Learning in conjunction with ICML, 2008.
- James Dickey. Scientific reporting and personal probabilities: Student's hypothesis. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 285–305, 1973.
- Thomas G. Dietterich. Approximate statistical tests for comparing supervised classification learning algorithms. *Neural Computation*, 10:1895–1924, 1998.
- Ward Edwards, Harold Lindman, and Leonard J Savage. Bayesian statistical inference for psychological research. *Psychological review*, 70(3):193, 1963.
- Andrew Gelman. Prior distributions for variance parameters in hierarchical models (comment on article by browne and draper). *Bayesian analysis*, 1(3):515–534, 2006.
- Andrew Gelman, Jennifer Hill, and Masanao Yajima. Why we (usually) don't have to worry about multiple comparisons. *Journal of Research on Educational Effectiveness*, 5 (2):189–211, 2012.

- Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian Data Analysis*. CRC press, 2013.
- Andrew Gelman, John B Carlin, Hal S Stern, and Donald B Rubin. *Bayesian Data Analysis*, volume 2. Taylor & Francis, 2014.
- Myles Hollander, Douglas A Wolfe, and Eric Chicken. *Nonparametric Statistical Methods*, volume 751. John Wiley & Sons, 2013.
- Miguel A Juárez and Mark FJ Steel. Model-based clustering of non-gaussian panel data based on skew-t distributions. *Journal of Business & Economic Statistics*, 28(1):52–66, 2010.
- Robert E Kass and Adrian E Raftery. Bayes factors. *Journal of the American Statistical Association*, 90(430):773–795, 1995.
- John K. Kruschke. Bayesian data analysis. Wiley Interdisciplinary Reviews: Cognitive Science, 1(5):658–676, 2010.
- John K. Kruschke. Bayesian estimation supersedes the t-test. *Journal of Experimental Psychology: General*, 142(2):573, 2013.
- John K. Kruschke. Doing Bayesian Data Analysis: A Tutorial with R, Jags and Stan. Academic Press, 2015.
- John K. Kruschke and Torrin M Liddell. The Bayesian New Statistics: Two Historical Trends Converge. *Available at SSRN 2606016*, 2015.
- Alexandre Lacoste, François Laviolette, and Mario Marchand. Bayesian comparison of machine learning algorithms on single and multiple datasets. In *Proc. of the Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS-12)*, pages 665–675, 2012.
- Bruno Lecoutre and Jacques Poitevineau. The Significance Test Controversy Revisited. Springer, 2014.
- Francesca Mangili, Alessio Benavoli, Cassio P. de Campos, and Marco Zaffalon. Reliable survival analysis based on the Dirichlet Process. *Biometrical Journal*, 57:10021019, 2015.
- Kevin P Murphy. Machine Learning: a Probabilistic Perspective. MIT press, 2012.
- Claude Nadeau and Yoshua Bengio. Inference for the generalization error. *Machine Learning*, 52(3):239–281, 2003.
- Steven L. Salzberg. On comparing classifiers: Pitfalls to avoid and a recommended approach. Data Mining and Knowledge Discovery, 1:317–328, 1997.
- Zbynek Sidak, Pranab Sen, and Jaroslav Hajek. *Theory of Rank Tests*. Probability and Mathematical Statistics. Elsevier Science, 1999.

- David Trafimow and Michael Marks. Editorial. Basic and applied social psychology, 37(1): 1–2, 2015.
- Peter Walley. Inferences from multinomial data: learning about a bag of marbles. *Journal* of the Royal Statistical Society. Series B (Methodological), 58(1):3–57, 1996.
- Ronald L Wasserstein and Nicole A Lazar. The ASA's statement on p-values: context, process, and purpose. *The American Statistician*, (just-accepted):00–00, 2016.
- Peter H Westfall, S Stanley Young, and S Paul Wright. On adjusting p-values for multiplicity. *Biometrics*, 49(3):941–945, 1993.
- Ian H Witten, Eibe Frank, and Mark Hall. Data Mining: Practical Machine Learning Tools and Techniques (third edition). Morgan Kaufmann, 2011.