

Statistical comparison of classifiers through Bayesian hierarchical modelling

Giorgio Corani · Alessio Benavoli · Janez Demšar · Francesca Mangili · Marco Zaffalon

Received: date / Accepted: date

Abstract We propose a new approach for the statistical comparison of algorithms which have been cross-validated on multiple data sets. It is a Bayesian hierarchical method; it draws inferences on single and on multiple datasets taking into account the mean and the variability of the cross-validation results. It is able to detect equivalent classifiers and to claim significances which have a practical impact. On each data sets it estimates more accurately than the existing methods the difference of accuracy between the two classifiers thanks to shrinkage. Such advantages are demonstrated by simulations on synthetic and real data.

1 Introduction

The statistical comparison of two competing algorithms is fundamental in machine learning. We take as an example throughout this paper the comparison of the accuracy of two classifiers. Typically the two classifiers are assessed via cross-validation and then compared via hypothesis testing.

Let us introduce some notation. We have a collection of q data sets; the actual mean difference of accuracy on the i -th data set is δ_i . On the i -th data set we obtain by cross-validation the measures $x_{i1}, x_{i2}, \dots, x_{in}$; they are cross-correlated with correlation ρ because of the overlapping training sets built during cross-validation. Their sample mean and standard deviation are \bar{x}_i and s_i . The maximum likelihood estimator (MLE) of δ_i is \bar{x}_i . The average difference of accuracy on the population of data sets (of which we observed only q data sets) is δ_0 .

The recommended approaches for comparing two classifiers on a single and on multiple data set are the correlated t-test and the signed-rank test (Demšar, 2006).

G. Corani, A. Benavoli F. Mangili and M. Zaffalon are with
Istituto Dalle Molle di studi sull'Intelligenza Artificiale (IDSIA), Manno, Switzerland
J. Demšar is with
Faculty of Computer and Information Science, University of Ljubljana, Slovenia
E-mail:
giorgio{alessio,francesca,zaffalon}@idsia.ch
janez.demsar@fri.uni-lj.si

The correlated t-test (Nadeau & Bengio, 2003) analyzes the cross-validation result on a single data set. It makes inference on about δ_i by analyzing the sampling distribution of a t statistic, corrected for correlation. On the i -th data set such statistic is $t = \bar{x}_i / \sqrt{\hat{s}_i^2 (\frac{1}{n} + \frac{\rho}{1-\rho})}$. The denominator of the statistic is the standard error; it is informative about the accuracy of \bar{x}_i as an estimator of δ_i . The standard errors largely varies across data sets, as a result of each data set having its own size and complexity.

The signed-rank test is applied after cross-validation. It makes inference about δ_0 analyzing the \bar{x}_i 's but ignoring the standard errors. It simplistically assumes the \bar{x}_i 's to be i.i.d., overlooking important pieces of information. There is no nhst test able to make inference about δ_0 accounting for the variability of cross-validation estimates.

Both the signed-rank and the correlated t-tests suffer from the shortcomings which characterize the null-hypothesis significance tests (nhst). For instance, the claimed statistical significances do not necessarily imply practical significance. A nhst rejects the null hypothesis when the p-value is smaller than the test size α ; yet the p-value depends (Wagenmakers, 2007) both on the effect size (the actual difference between the two classifiers) and the sample size (the number of collected observations). Thus you can easily reject the null hypothesis of the signed-rank test, even if the two classifier are almost equivalent; to this end it is enough to compare them on a large enough collection of data sets. Null hypotheses can virtually always be rejected with enough data (Lecoutre & Poitevineau, 2014, Chap.5.2).

Moreover the nhst cannot verify the null hypothesis; it can only *reject* it or *fail to reject* it (Kruschke, 2015, Chap.11). When the nhst fails to reject the null hypothesis, it does *not* conclude that the two classifiers are equivalent. Instead it draws a non-committal conclusion: there is not enough evidence for rejecting the null hypothesis, which might be true or not.

The problem can be solved switching to the Bayesian approach and setting a region of practical equivalence (rope) (Kruschke, 2013). However there is currently no Bayesian approach able to make inference on both the δ_i 's and δ_0 .

The Bayesian correlated t-test (Corani & Benavoli, 2015) computes the posterior distribution of δ_i on a *single* data set. Lacoste et al. (2012) compares two classifiers on multiple data sets by modelling each data set as an independent Bernoulli trial. Its possible outcomes are the first classifier being more accurate than the second or vice versa. Each Bernoulli trial is characterized by different probabilities. Thus the data sets are assumed to be independent and not identically distributed (i.-i.d.), trying to overcome the i.i.d. assumption. One eventually computes the probability of the first classifier being more accurate than the second classifier on more than half of the q data sets. Yet the achieved conclusion applies only to the q available data sets; it does not generalize to the population of data sets and so there is no inference about δ_0 .

We propose a novel model which overcomes such limits. It is the first model which represent the two stochastic components of the process of cross-validation on multiple data sets: a) the distribution of the difference of accuracy across the population of data sets (high-level distribution); b) the distribution of the cross-validation measures on the i -th data set given δ_i .

It models the fact that the δ_i 's are drawn from a high-level distribution with mean δ_0 . It also assumes the cross-validation measures on the i -th data to be characterized by their own standard deviation σ_i 's, also drawn from a high-level distribution. This is a hierarchical Bayesian model since it has multiple levels of unknown parameters.

We want to make inference about the δ_i 's and δ_0 . By applying the rope on the posterior distribution of the δ_i 's and the δ_0 we are able to detect equivalent classifiers and claim significance which have a practical impact. We adopt as rope the interval $(-0.01, 0.01)$. We claim two classifiers to be practically equivalent when a large part of the posterior probability is concentrated within the rope. Conversely we declare two classifiers as significantly different if a large part of the posterior probability is concentrated at the left (or at the right) of the rope.

A further merit of the hierarchical model is that it jointly estimates the δ_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 shrinkage estimator is more accurate than MLE in the case of uncorrelated data; see the discussion in (Murphy, 2012, Sec.6.3.3.2). Yet there are currently no results about shrinkage estimation with correlated data, such as those yielded by cross-validation. We prove that also in the correlated case shrinkage outperforms MLE. Our 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. The hierarchical model thus estimates the δ_i 's more accurately than the existing tests.

2 Existing approaches

Throughout this paper we take accuracy as an example of indicator of performance. However our discussion readily applies to any other indicator of performance. Assume that we have performed m runs of k -folds cross-validation on each data set, providing both classifiers with the same training and test sets. The *differences of accuracy* on each fold of cross-validation are $\mathbf{x}_i = \{x_{i1}, x_{i2}, \dots, x_{in}\}$, where $n = mk$. The x_{ij} 's ($j=1,2,\dots,n$) are correlated because of the overlapping training sets adopted by cross-validation. Nadeau & Bengio (2003) prove that there is no unbiased estimator of the correlation and they approximate it as $\rho = \frac{1}{k}$, where k is the number of folds. They devise the correlated t -test, which corrects the standard error of the t -test by accounting for the correlation. This is the standard approach for comparing two classifiers on a single data set. The signed-rank test is instead the recommended method (Demšar, 2006) to compare two classifiers on a collection of q different data sets, after having performed cross-validation on each data set. The test analyzes the vector $\bar{\mathbf{x}}$ constituted by the mean differences measured on each data set: $\bar{\mathbf{x}} = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_q\}$, assuming them to be i.i.d.. Both the correlated t -test and the signed-rank are *nhst* tests. There are also Bayesian approaches for comparing two competing classifiers.

The Bayesian correlated t -test (Corani & Benavoli, 2015) adopts a generative model which jointly generates n observations from random variables which are jointly normal, equally cross-correlated and which have the same variance. It adopts the same correlation heuristic of Nadeau & Bengio (2003). Under a conju-

gate non-informative prior, the posterior distribution of δ_i matches the sampling distribution of the frequentist correlated t-test. Thus the posterior expected value of δ_i is \bar{x}_i .

Lacoste et al. (2012) compare two classifiers on multiple data sets by modelling each data set as an independent Bernoulli trial, whose possible outcomes are either the first classifier being more accurate than the second or vice versa. The probabilities of such two outcomes are computed by another model estimated independently in each data set, but without managing the correlation of cross-validation results. Thus the data sets are assumed to be i.-i.d.. The number of data set in which the first classifier is more accurate than the second follows a Poisson-binomial distribution. This approach makes no inference about δ_0 and does not tries to estimate how the two classifiers compare on the population of data sets. However it has been followed also by (Corani & Benavoli, 2015), coupling it with the posterior probabilities yielded by the Bayesian correlated t-test.

3 The hierarchical model

We propose a new method for comparing two classifiers. It draws inferences about the population of datasets taking into consideration both mean and variability of the cross-validation results on the individual data sets. It estimates the individual δ_i 's more accurately than the MLE. It makes inferences on multiple datasets by estimating the ground truth difference between the classifiers.

The method we propose is a Bayesian hierarchical model. Its main assumptions are described by the following probabilistic model:

$$\mathbf{x}_i \sim MVN(\mathbf{1}\delta_i, \mathbf{\Sigma}_i) \quad (1)$$

$$\delta_1 \dots \delta_q \sim t(\delta_0, \sigma_0, \nu) \quad (2)$$

$$\sigma_1 \dots \sigma_q \sim \text{unif}(0, \bar{\sigma}) \quad (3)$$

Equation (1) 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 generated from random variables which are equally cross-correlated (ρ), which have the same mean (δ_i) and variance (σ_i) and which are jointly normal. In fact, it states that, for each dataset i , \mathbf{x}_i is multivariate normal with mean $\mathbf{1}\delta_i$ (where $\mathbf{1}$ is a vector of ones) and covariance matrix $\mathbf{\Sigma}_i$ defined as follows: the diagonal elements are σ_i^2 and the out-of-diagonal elements are $\rho\sigma_i^2$, where $\rho = \frac{n_{te}}{n_{tr}}$. This idea is borrowed from Corani & Benavoli (2015). The normality assumptions is sound since the average of an indicator over the instances of the test set tends to be normally distributed by the central limit theorem. Equation (2) models the fact that the mean difference of accuracies in the single datasets, δ_i , depends on δ_0 that is the ‘‘ground truth’’ difference between the classifiers.

We assume the δ_i 's to be drawn from a high-level Student distribution with mean δ_0 , variance σ_0^2 and degrees of freedom ν . The Student distribution robustly deals with outliers thanks to its heavy tails (Gelman et al., 2014; Kruschke, 2013). Thus the hierarchical model can robustly estimate δ_0 even in presence of some δ_i 's located far away from the others. Moreover the Student distribution is more expressive than the Gaussian and thus can it generally fits better the data.

The hierarchical model assigns to the i-th data set its own standard deviation σ_i , assuming the σ_i 's to be drawn from a common distribution, see Eqn.(3). 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 σ_i 's is $\text{unif}(0, \bar{\sigma})$, as recommended by Gelman (2006), as it yields inferences which are insensitive on $\bar{\sigma}$ if $\bar{\sigma}$ is large enough. To this end we set $\bar{\sigma} = 1000 \cdot \bar{s}$ (Kruschke, 2013) where $\bar{s} = \sum_i^q s_i/q$.

We refine the model with the prior on some further parameters. The height of the tails of the Student distribution is controlled by the degrees of freedom ν . When ν is small, the distribution has heavy tails; when ν is large (e.g., above 30), the distribution is nearly normal. Two different priors for ν are proposed in literature: the shifted exponential (Kruschke, 2013) and the Gamma(2,0.1) (Juárez & Steel, 2010). We reparametrize the shifted exponential of Kruschke (2013) as a Gamma(1,0.0345). We found the inferences of the hierarchical model to be sensitive on the choice of $p(\nu)$. We address the problem by adopting a hierarchical approach, thus assuming $p(\nu) = \text{Gamma}(\alpha, \beta)$, with $\alpha \sim \text{unif}(\underline{\alpha}, \bar{\alpha})$ and $\beta \sim \text{unif}(\underline{\beta}, \bar{\beta})$, setting $\underline{\alpha}=1$, $\bar{\alpha}=2$, $\underline{\beta}=0.01$, $\bar{\beta}=0.1$. The resulting model is robust to perturbations of the priors on α and β .

We still need a prior for the scale parameter σ_0 and the location parameter δ_0 of the Student. We set $\sigma_0 \sim \text{unif}(0, \bar{\sigma}_0)$, with $\bar{\sigma}_0 = 1000 \cdot s_{\bar{x}}$, where $s_{\bar{x}}$ is the standard deviation of the \bar{x}_i 's. We adopt an improper uniform for δ_0 .

These consideration are reflected by the following probabilistic model:

$$\nu \sim Ga(\alpha, \beta) \quad (4)$$

$$\alpha \sim \text{unif}(\underline{\alpha}, \bar{\alpha}) \quad (5)$$

$$\beta \sim \text{unif}(\underline{\beta}, \bar{\beta}) \quad (6)$$

$$\delta_0 \sim \text{unif}(-1, 1) \quad (7)$$

$$\sigma_0 \sim \text{unif}(0, \bar{\sigma}_0) \quad (8)$$

The limits of the uniform distribution of $p(\delta_0)$ of Eqn.(8) are suitable to deal with indicators whose difference bounded between 1 and -1, such as accuracy, AUC, precision and recall. In order to deal with an unbounded indicator such as log-loss, one should instead adopt an improper prior as $p(\delta_0)$.

We implemented the hierarchical model in Stan (<http://mc-stan.org>) (Carpenter et al., 2016), a language for Bayesian inference. We vectorized the computation of the likelihood over the q data sets, reducing the computational time of about one order of magnitude compared to the non-vectorized version. The analysis of the results of 10 runs of 10-folds cross-validation on 50 data sets (that means a total of 5000 observations) takes about three minutes on a standard computer. We make available the Stan code of the hierarchical model from <https://github.com/BayesianTestsML/>, within the repository `tutorial/hierarchical`. The same repository provides the R code of the simulations of Sec. 5.

3.1 The inference of the test

Once the hierarchical model has been learned, it can be queried in different ways. You can make inference on the difference of accuracy on the i -th data set by inspecting the posterior distribution of δ_i . To make inference on the average difference between the two classifiers on the population of data sets, you have to check the posterior distribution of the δ_0 .

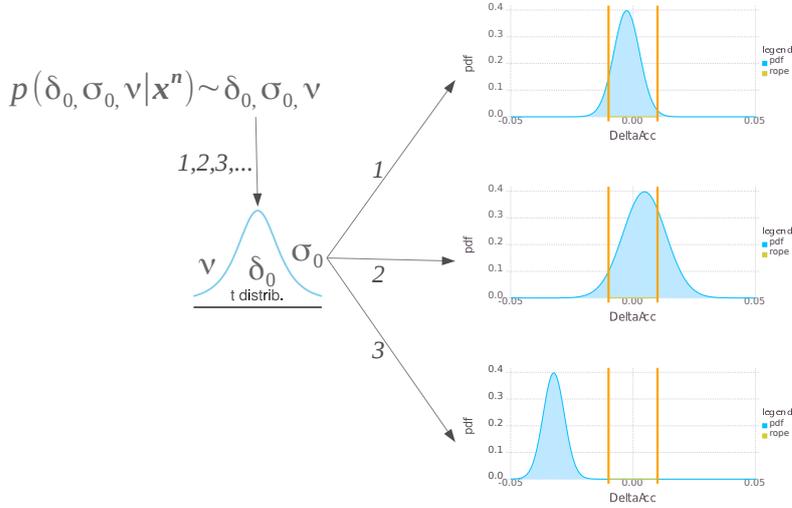


Fig. 1 Sampling and posterior inference schema

We instead focus on estimating the posterior distribution of the difference of accuracy between the two classifiers on a *future unseen data set*. For comparing classifiers, this is the most important inference we are interested in.

To compute such inference, we proceed as follows:

1. initialize the counters $n_{left} = n_{rope} = n_{right} = 0$;
2. for $i = 1, 2, 3, \dots, N_s$ repeat
 - sample μ_0, σ_0, ν from the posterior of these parameters;
 - define the posterior of the mean difference accuracy on the next dataset, i.e., $t(\delta_{next}; \delta_0, \sigma_0, \nu)$;
 - from $t(\delta_{next}; \delta_0, \sigma_0, \nu)$ compute the three probabilities $p(left)$ (integral on $(-\infty, r]$), $p(rope)$ (integral on $[-r, r]$) and $p(right)$ (integral on $[r, \infty)$);
 - determine the highest among $p(left), p(rope), p(right)$ and increment the respective counter $n_{left}, n_{rope}, n_{right}$;
3. compute $P(left) = n_{left}/N_s$, $P(rope) = n_{rope}/N_s$ and $P(right) = n_{right}/N_s$;
4. decision: when $P(rope) > 1 - \alpha$ declare the two classifiers to be *practically equivalent*; when $P(left) > 1 - \alpha$ or $P(right) > 1 - \alpha$ we declare the two classifiers to be significantly different in the respective directions.

We have chosen $r = 0.01$ and, thus, our region of practical equivalence (rope) is $[-0.01, 0.01]$. Figure 1 shows a diagram of this inference schema and reports three sampled posteriors $t(\delta_{next}; \delta_0, \sigma_0, \nu)$. For these three cases we have that $(p(left), p(rope), p(right))$ are respectively (from top to bottom) $(0.08, 0.90, 0.02)$, $(0.05, 0.67, 0.28)$, $(1, 0, 0)$ and so after these three steps $n_{left} = 1$, $n_{rope} = 2$, $n_{right} = 0$ (in the next experiments we will consider $N_s = 4'000$).

We can further explain this procedure through an example. Consider a bag of special coins (our posterior samples) which can land tails, head or remain standing

(equiv. to left, right or rope). We want to make inference on the type of bias (left, right or rope) of the population of coins. Then, we extract a coin and determine its type (i.e., we determine the highest value between $p(left), p(right), p(rope)$). We repeat this procedure many times (N_s) determining the distribution of the type of bias of the coins inside the bag ($P(left), P(rope), P(right)$). If for instance $P(right)$ exceeds $1-\alpha$ (e.g., 0.95) then we decide that the bag of coins is of type *right* (biased towards Head). This means that there is a difference between the classifiers (with high probability the population of coins is biased towards a particular direction). A similar inference is performed by the Bayesian signed-rank test (Benavoli et al., 2014).

This way of taking automatic decisions is suitable to perform automatic experiments as we work with synthetic data in Section 5.2. When however one presents a novel classifier on a real case study we suggest to *reason* about the obtained posterior probabilities rather than taking automatic decisions. The posterior distribution provides much more information than its integrals over the three regions left, right and rope; for instance in Figure 1 the second posterior is wider (more uncertainty) than the other two cases.

4 The shrinkage estimator for cross-validation

The hierarchical model jointly estimates the δ_i 's by applying shrinkage to the \bar{x}_i 's. In the uncorrelated case, the shrinkage estimator is known to be more accurate than the MLE. In this section we show that the shrinkage estimator is more accurate than MLE also in the correlated case, such as the data generated by cross-validation. This allows the hierarchical model to be more accurate than the existing method in the estimation of the δ_i 's.

The δ_i 's of the hierarchical model are independent given the parameters of the higher-level distribution. If such parameters were known, the δ_i 's would be conditionally independent and they would be independently estimated. Instead such parameters are unknown, causing the δ_0 and the δ_i 's to be jointly estimated. As a result the estimate of each δ_i is informed by data collected also on all the other data sets. Intuitively, each data set informs the higher-level parameters, which in turn constrain and improves the parameters of the individual data sets (Kruschke, 2013, Chap.9).

To show this, we assume the cross-validation results on the q data sets to be generated by the hierarchical model:

$$\begin{aligned}\delta_i &\sim p(\delta_i) \\ \mathbf{x}_i &\sim MVN(\mathbf{1}\delta_i, \Sigma)\end{aligned}\tag{9}$$

where, to simplify the analytical analysis, we have assumed the variances σ_i^2 of the individual data sets to be equal to σ^2 and known. Thus all data sets have the same covariance matrix Σ , which is defined as follows: the variances are all σ^2 and the correlations are all ρ . Note that Eqn.(9) coincides with (1). This is a general model which makes no assumptions about the distribution $p(\delta_i)$. We denote the two first moments of $p(\delta_i)$ as $E[\delta_i] = \delta_0$ and $\text{Var}[\delta_i] = \sigma_0^2$.

For the purpose of analytic analysis we study the MAP estimates of the parameters $\delta_1, \dots, \delta_m, \delta_0, \sigma_0^2$, which asymptotically tend to the Bayesian estimates. A

hierarchical model is being fitted to the data. Such model is a simplified version of that presented in Sec.3. In particular $p(\delta_i)$ is Gaussian for analytical tractability.

$$P(\bar{\mathbf{x}}, \boldsymbol{\delta}, \delta_o, \sigma_o^2) = \prod_{i=1}^q N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma})N(\delta_i; \delta_o, \sigma_o^2)p(\delta_o, \sigma_o^2) \quad (10)$$

This model is misspecified since $p(\delta_i)$ is generally not Gaussian. Nevertheless, it correctly estimates the mean and variance of $p(\delta_i)$, as we show in the following.

Proposition 41 *The derivatives of the logarithm of $P(\bar{\mathbf{x}}, \boldsymbol{\delta}, \delta_o, \sigma_o^2)$ are:*

$$\begin{aligned} \frac{d}{d\delta_i} \ln(P(\cdot)) &= \frac{\delta_o - \delta_i}{\sigma_o^2} + \frac{\bar{x}_i - \delta_i}{\sigma_n^2} \\ \frac{d}{d\delta_o} \ln(P(\cdot)) &= \frac{-q\delta_o + \sum_{i=1}^q \delta_i}{\sigma_o^2} + \frac{d}{d\delta_o} \ln(p(\delta_o, \sigma_o^2)) \\ \frac{d}{d\sigma_o} \ln(P(\cdot)) &= \frac{q\delta_o^2 + \sum_{i=1}^q \delta_i^2 - 2\delta_o \sum_{i=1}^q \delta_i - q\sigma_o^2}{\sigma_o^3} + \frac{d}{d\sigma_o} \ln(p(\delta_o, \sigma_o^2)) \end{aligned}$$

If we further assume that $p(\delta_o, \sigma_o^2) \approx \text{constant}$ (flat prior), by equating the derivatives to zero, we derive the following consistent estimators:

$$\delta_o = \frac{\sum_{i=1}^q \hat{\delta}_i}{q}, \quad \sigma_o^2 = \frac{1}{q} \sum_{i=1}^q (\hat{\delta}_i - \hat{\delta}_o)^2 \quad (11)$$

$$\hat{\delta}_i = \frac{\hat{\sigma}_o^2 \bar{x}_i + \sigma_n^2 \frac{1}{q} \sum_{i=1}^q \bar{x}_i}{\hat{\sigma}_o^2 + \sigma_n^2} = w\bar{x}_i + (1-w)\frac{1}{q} \sum_{i=1}^q \bar{x}_i. \quad (12)$$

where $w = \hat{\sigma}_o^2 / (\hat{\sigma}_o^2 + \sigma_n^2)$ and, to keep a simple notation, we have not explicitated the expression $\hat{\sigma}_o$ as a function of \bar{x}_i, σ_n^2 . Notice that the estimator $\hat{\delta}_i$ shrinks the estimate towards $\frac{1}{q} \sum_{i=1}^q \bar{x}_i$ that is an estimate of δ_o . Hence, the Bayesian hierarchical model consistently estimates δ_o and σ_o^2 from data and converges to the shrinkage estimator $\hat{\delta}_i(\mathbf{x}_i) = w\bar{x}_i + (1-w)\delta_o$.

It is known that the shrinkage estimator achieves a lower error than MLE in case of uncorrelated data; see (Murphy, 2012, Sec.6.3.3.2) and the references therein. However there is currently no analysis of shrinkage with correlated data, such as those yielded by cross-validation. We study this problem in the following.

Consider the generative model (9). The likelihood regarding the i -th data set is:

$$p(\mathbf{x}_i | \delta_i, \boldsymbol{\Sigma}) = N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma}) = \frac{\exp(-\frac{1}{2}(\mathbf{x}_i - \mathbf{1}\delta_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \mathbf{1}\delta_i))}{(2\pi)^{n/2} \sqrt{|\boldsymbol{\Sigma}|}}. \quad (13)$$

Let us denote by $\boldsymbol{\delta}$ the vector of the δ_i 's. The joint probability of data and parameters is:

$$P(\boldsymbol{\delta}, \mathbf{x}_1, \dots, \mathbf{x}_q) = \prod_{i=1}^q N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma})p(\delta_i)$$

Let us focus on the i -th group, denoting by $\hat{\delta}_i(\mathbf{x}_i)$ an estimator of δ_i . The mean squared error (MSE) of the estimator w.r.t. the true joint model $P(\delta_i, \mathbf{x}_i)$ is:

$$\iint (\delta_i - \hat{\delta}_i(\mathbf{x}_i))^2 N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma})p(\delta_i)d\mathbf{x}_i d\delta_i. \quad (14)$$

Proposition 42 *The MSE of the maximum likelihood estimator is:*

$$\begin{aligned} \text{MSE}_{\text{MLE}} &= \iint (\delta_i - \bar{x}_i)^2 N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma})p(\delta_i)d\mathbf{x}_i d\delta_i \\ &= \frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1}, \end{aligned}$$

which we denote in the following also as $\sigma_n^2 = \frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1}$.

Now consider the shrinkage estimator $\hat{\delta}_i(\mathbf{x}_i) = w\bar{x}_i + (1-w)\delta_0$ with $w \in (0, 1)$, which pulls the MLE estimate \bar{x}_i towards the mean δ_0 of the upper-level distribution.

Proposition 43 *The MSE of the shrinkage estimator is:*

$$\begin{aligned} \text{MSE}_{\text{SHR}} &= \iint (\delta_i - w\bar{x}_i - (1-w)\delta_0)^2 N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma})p(\delta_i)d\mathbf{x}_i d\delta_i \\ &= w^2 \sigma_n^2 + (1-w)^2 \sigma_0^2. \end{aligned}$$

As we have seen, the hierarchical model converges to the shrinkage estimator with $w = \sigma_0^2 / (\sigma_0^2 + \sigma_n^2)$. Then:

$$\begin{aligned} \text{MSE}_{\text{SHR}} &= w^2 \sigma_n^2 + (1-w)^2 \sigma_0^2 = \frac{\sigma_0^4 + \sigma_n^2 \sigma_0^2}{(\sigma_0^2 + \sigma_n^2)^2} \sigma_n^2 \\ &= \frac{\sigma_0^2}{(\sigma_0^2 + \sigma_n^2)} \sigma_n^2 < \sigma_n^2 = \text{MSE}_{\text{MLE}}. \end{aligned}$$

Therefore, the shrinkage estimator achieves a smaller mean squared error than the MLE.

5 Experiments

Previous studies (Bouckaert, 2003; Kohavi, 1995) recommend to perform 10 runs of 10-folds cross-validation and we follow this practice in our experiments. We have thus $n=100$ observations on each data set.

We start by presenting the results on synthetic data. The synthetic generation of data proceeds as follows. We assume a high-level distribution with mean δ_0 , from which we sample the δ_i 's. We adopt different high-level distributions depending on the experiment being carried out. Let us denote by $p(\delta_i)$ the actual distribution from which we sample the δ_i 's and by $\hat{p}(\delta_i)$ the high-level Student distribution fitted by the hierarchical model.

We then implement on each data set the cross-validation of two classifiers. On the i -th data set we simulate two classifier whose actual mean difference of accuracy is δ_i , following the method of Corani & Benavoli (2015). The simulation returns the cross-validation measures $x_{i1}, x_{i2}, \dots, x_{i100}$, whose mean is \bar{x}_i . We analyze such results using the hierarchical model of Section 3 and the other existing methods.

5.1 Improved inferences on the individual data sets

In this experiment we assess the inferences on the individual data sets. We choose a mixture as $p(\delta_i)$, thus making the Student distribution of the hierarchical model misspecified. We adopt the strongly bimodal mixture $p(\delta_i) = \sum_{i=1}^k \pi_k N(\delta_i | \mu_k, \sigma_k)$ with $k=2$, $\mu_1=0.005$, $\mu_2=0.02$, $\sigma_1=\sigma_2=\sigma=0.001$, $\pi_1 = \pi_2 = 0.5$.

We consider the following number of data sets: $q = \{10, 20, 30, 40, 50\}$. For each value of q we repeat 500 experiments organized as follows: sampling of the δ_i 's from the mixture; cross-validation of the two competing classifiers on each data set; inference of the hierarchical model. Moreover we perform on each data set the Bayesian correlated t-test. We then measure MSE_{MLE} and MSE_{SHR} . The shrinkage estimator has no closed form due to the complexity of the model of Section 3, and we thus compute it numerically.

As shown in Tab.1, MSE_{SHR} is considerably lower than MSE_{MLE} . The shrinkage estimator reduces mse as q increases: this happens because more data sets allow to estimate more accurately the parameters of the high-level distribution, improving the shrinkage. This learning effect cannot be replicated by the existing methods which estimate independently the parameters of each data set.

q	Mean Squared Error	
	MLE	Shrinkage
Mixture experiment		
5	.00036	.00017
10	.00036	.00014
50	.00036	.00012
Gaussian experiment		
5	.00036	.00020
10	.00036	.00014
50	.00036	.00012

Table 1 Inferences regarding individual data sets. The scale of the actual errors on the estimation of the δ_i 's can be realized considering that for instance $0.02^2=.0004$.

For $q=50$ (a common size for a machine learning study), the hierarchical model reduces MSE of about 60% compared to the maximum likelihood estimator (Tab 1). This happens *despite* the severe mismatch between $p(\delta_i)$ and $\hat{p}(\delta_i)$, confirming the formal analysis of the previous section: the shrinkage estimator dominates the MLE under the mild assumption that $\hat{p}(\delta_i)$ reliably estimates the first two moments of $p(\delta_i)$.

As a double-check, we repeated the same experiments using a Gaussian distribution $p(\delta_i)$, with the same mean and variance of the mixture. The results are consistent among the two experiments (Tab 1), further proving the correctness of the analytical analysis. This hierarchical model for the first time jointly estimates the parameter referring to different data sets; for this reason it delivers the best estimates of the δ_i 's so far.

5.2 Simulation of formally equivalent classifiers

From now on we adopt a Cauchy distribution as a more realistic model for $p(\delta_i)$. In the following experiment we simulate the null hypothesis of the signed-rank, setting $\delta_0 = 0$. We set the scale factor of the distribution to 1/6 of the rope length. We consider the following number of data sets: $q = \{10, 20, 30, 40, 50\}$. For each value of q we repeat 500 experiments as in the previous section.

The signed-rank test proves to be correctly calibrated: it rejects the null hypothesis about 5% of the times regardless the sample size. This is its expected behavior when H_0 is true. Notice that this behavior, despite being formally correct, provides no valuable insight. When the signed-rank fails to reject the null (95% of the times), it takes a non committal conclusion: it is *not* claiming that the null hypothesis is true. It simply states that there is not enough evidence for rejecting it. When the signed-rank rejects the null (5% of the times), it draws instead a wrong conclusion since $\delta_0=0$.

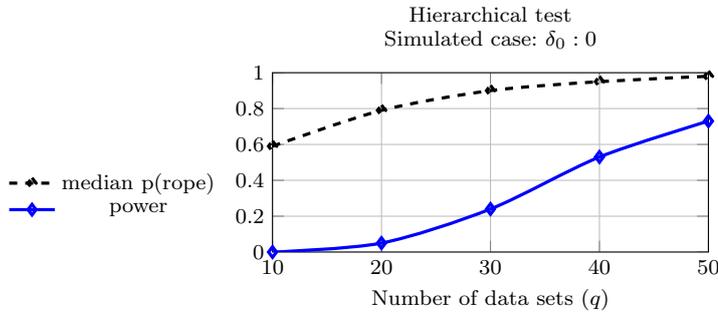


Fig. 2 Behavior of the hierarchical classifier when simulating the null hypothesis of the signed-rank test ($\delta_0 = 0$). The power is the proportion of simulations in which the hierarchical test estimates $p(\text{rope}) > 0.95$.

The hierarchical model shows a more sensible behavior. The posterior probability of rope steadily increases with q : when there are more data sets, there is more information for estimating the difference between the two classifiers. For $q=50$ (the typical size of a machine learning study), the posterior probability of rope is on average well above 90% (Fig.2). A second characterization of the behavior of the test is provided by its power, namely the proportion of simulations in which the test estimates $p(\text{rope}) > 0.95$. The power of the hierarchical test increases with q , reaching about 0.7 for $q=50$.

Another strength of the hierarchical test is that it makes no Type I error: in our simulations it never estimates $p(\text{left}) > 95\%$ or $p(\text{right}) > 95\%$. It is indeed known that Bayesian estimation with rope drastically reduces the Type I errors (Kruschke, 2013) compared to nhst.

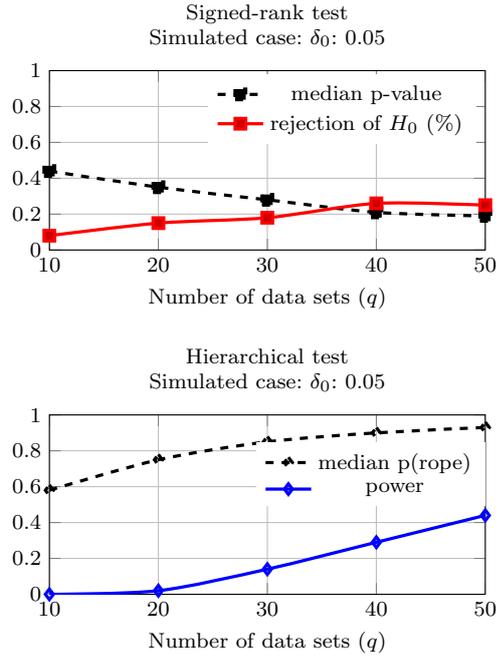


Fig. 3 Simulation of two practically equivalent ($\delta_0 = 0.005$) classifiers. Upper plot: the p-value of the signed-rank decreases with q and thus the test rejects more often H_0 as q increases. Lower plot: the hierarchical test is able to recognize that the two classifiers are practically equivalent. The probability of rope increases with q and so does the power (proportion of simulations in which the model estimates $p(\text{rope}) > 0.95$).

5.3 Simulation of practically equivalent classifiers

In these experiments we simulate two classifiers whose actual difference of accuracy is small but different from zero. This is a common situation in practice. We set $\delta_0 = 0.005$, a difference which has no practical value according to our definition of rope. We consider $q = \{10, 20, 30, 40, 50\}$. For each value of q we repeat 500 experiments. The power (% of rejections of H_0) of the signed-rank increases with q (Fig.3). As already discussed, you can reject the null of the signed rank when comparing two almost equivalent classifiers: it is just a matter of comparing them on enough data sets. This is clearly shown by Fig.3. When 50 data sets are available, the signed-rank rejects the null in about 25% of the simulations, despite the trivial value of δ_0 .

The hierarchical behaves robustly. It is only slightly less powerful in recognizing equivalence (Fig.3) than in the previous experiment since δ_0 is now closer to the limit of the rope. It responds to the increase of q by increasing the posterior probability of rope. Also in this condition it makes no Type I errors.

5.4 Simulation of practically different classifiers

We now simulate two classifiers which are significantly different. We consider different values of δ_0 : $\{0.01, 0.02, 0.03\}$. We set the scale factor of the Cauchy to $\sigma_0=0.01$. We study how the power of the tests varies with δ_0 , fixing the number of data sets to $q=50$, the typical size of a machine learning study. We repeat 500 experiments for each value of q . The results are shown in Fig.4. The signed-rank test is more powerful, especially when δ_0 list just slightly outside the rope.

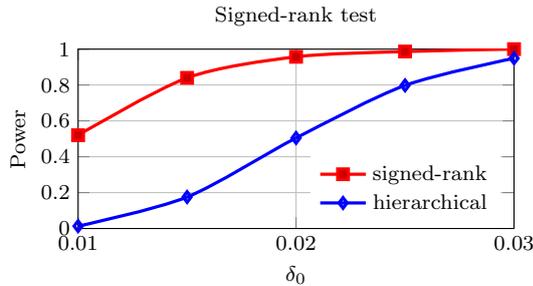


Fig. 4 Recognition of two significantly different classifiers. The power of signed-rank test is computed as the proportion of simulations in which it returns a p-value <0.05 . The power of the hierarchical model is the proportion of simulations in which it returns $p(\text{right}) > 0.95$.

Let us start with $\delta_0=0.01$. This is a borderline case: the mean of the distribution from which we sample the δ is exactly at the border between the rope and the right region. A correct answer is thus to estimate posterior probability of 50% for both regions. The hierarchical test provides sensible estimates: over the 500 experiments, the median $p(\text{rope})$ is 0.56 while the median $p(\text{right})$ is 0.44. The power of the hierarchical test then increases with δ_0 , though remaining lower than that of the signed-rank.

The power of a test is the percentage of rejections of the null hypothesis, in case in which it is false. For the hierarchical test, it is percentage of cases in which the probability of δ_0 belonging to the right region outside the rope exceeds 0.95. The two tests have about the same power for $\delta_0=0.03$ or higher.

6 Analysis on real data sets

We consider four classifiers: naive Bayes (nbc), hidden naive Bayes (hnb), decision tree (j48), grafted decision tree (j48gr). A description of all such classifiers can be found in Witten et al. (2011). We run all experiments using WEKA¹. We perform 10 runs of 10-folds cross-validation for each classifier on each data set, using 54 data sets from the WEKA data sets page. We compare the conclusions of the hierarchical model and of the signed-rank test. The results are given in Tab.2.

Both the signed-rank and hierarchical test identify hnb as significantly more accurate than naive Bayes; this is in agreement with the previous literature (Zhang et al., 2005).

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

		Signed rank		Hierarchical test		
<i>left</i>	<i>right</i>	<i>p value</i>	<i>p(left)</i>	<i>p(rope)</i>	<i>p(right)</i>	
nbc	hnb	0.00	0.00	0.00	1.00	
nbc	j48	0.46	0.20	0.01	0.79	
nbc	j48gr	0.39	0.15	0.01	0.84	
hnb	j48	0.07	0.91	0.07	0.03	
hnb	j48gr	0.08	0.92	0.05	0.03	
j48	j48gr	0.00	0.00	1.00	0.00	

Table 2 Comparison of real classifiers with the signed-rank test and the hierarchical test.

Moreover, both tests declare no significance when analyzing nbc vs j48 and nbc vs j48gr. hnb vs j48 and hnb vs j48gr. Yet the output of the hierarchical model is more informative. Take for instance nbc vs j48. A p-value of 0.46 conveys no information, as already discussed. The hierarchical model shows that there is negligible probability (0.01) of nbc and j48 being practically equivalent; a considerable probability of j48 being more accurate (0.79) and some probability of nbc being more accurate than j48 (0.20). When comparing j48 vs hnb, both test are close to declare hnb as significantly more accurate (with confidence 95%). A similar situation happens in the comparison between j48gr and hnb.

Interestingly, the two test draw opposite conclusions when comparing j48 and j48gr. The signed-rank declares j48gr to be significantly more accurate than j48 (p-value 0.00) while the hierarchical model declares them to be practically equivalent, with $p(\text{rope})=1$. The reason for this behavior can be understood by looking at the cross-validation results. As shown in Fig.5, the signs of the \bar{x}_i 's are mostly in favor of j48gr, causing the signed rank test to claim significance. Yet almost all the \bar{x}_i 's lie within the rope; thus the hierarchical model claims them to be practically equivalent. If the rope values are sensibly chosen, this looks indeed as an appropriate conclusion.

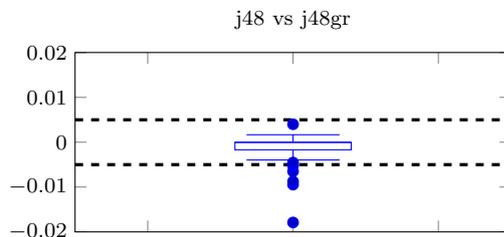


Fig. 5 Boxplots of the differences of accuracy \bar{x}_i 's between j48 and j48gr on 54 data sets. The dashed line shows the rope.

7 Sensitivity analysis

As a further control we inspected the posterior distribution of the variances σ_i 's. Dealing with correlated data, the maximum likelihood estimator of the variance

is not the usual one; it is instead the estimator given in the Appendix of (Corani & Benavoli, 2015). The posterior expected values of the variances consistently converge to the maximum likelihood estimator. The inferences of the model remain consistent if we adopt an empirical Bayes approach, substituting the prior (3) with a fixed value constituted by the maximum likelihood estimator.

8 Conclusions

The proposed hierarchical model provides a realistic model of the data generated by cross-validation across multiple data sets. It reliably detects practically equivalent classifiers and it claims statistical significances which have a practical meaning. On the individual data sets it yields more accurate inferences than the existing methods, being the first approach which jointly estimates the parameters referring to different data sets.

Acknowledgements

The research in this paper has been partially supported by the Swiss NSF grants ns. IZKSZ2_162188.

9 Appendix

9.1 Proofs

Proof of Proposition 41 Consider the hierarchical model:

$$\begin{aligned} P(\bar{\mathbf{x}}, \boldsymbol{\delta}, \delta_0, \sigma_0^2) \\ = \prod_{i=1}^q N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma})N(\delta_i; \delta_o, \sigma_o^2)p(\delta_o, \sigma_o^2) \end{aligned} \quad (15)$$

We aim to compute the derivative of the $\log(P(\bar{\mathbf{x}}, \boldsymbol{\delta}, \delta_0, \sigma_0^2))$ w.r.t. the parameter $\delta_i, \delta_o, \sigma_o^2$. Consider the quadratic term from the first and second Gaussian:

$$\frac{1}{2}(\mathbf{x}_i - \mathbf{1}\delta_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \mathbf{1}\delta_i) + \frac{1}{2\sigma_o^2}(\delta_i - \delta_o)^2$$

its derivatives w.r.t. δ_i is $\mathbf{1}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \mathbf{1}\delta_i) + \frac{1}{\sigma_o^2}(\delta_i - \delta_o)$. Exploiting the fact that

$$\begin{aligned} \mathbf{1}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \mathbf{1}\delta_i) &= \mathbf{1}^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \mathbf{1}\bar{x}_i + \mathbf{1}\bar{x}_i - \mathbf{1}\delta_i) \\ &= \mathbf{1}^T \boldsymbol{\Sigma}^{-1}(\mathbf{1}\bar{x}_i - \mathbf{1}\delta_i) \end{aligned}$$

it follows that

$$\frac{d}{d\delta_i} \ln(P(\cdot)) \propto \frac{1}{\sigma_n^2}(\bar{x}_i - \delta_i) + \frac{1}{2\sigma_o^2}(\delta_i - \delta_o)^2$$

where $\sigma_n^2 = \frac{1}{\mathbf{1}^T \boldsymbol{\Sigma}^{-1} \mathbf{1}} = \frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1}$. The latter equality can be derived by Corani & Benavoli (2015)[Appendix], i.e.,

$$\frac{1}{\mathbf{1}^T \boldsymbol{\Sigma}^{-1} \mathbf{1}} = \frac{n}{1 + (n-1)\rho} = \frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1}$$

The other derivatives can be computed easily.

Proof of Proposition 42 Let us consider the likelihood:

$$\begin{aligned} p(\mathbf{x}_i|\delta_i, \boldsymbol{\Sigma}) &= N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma}) \\ &= \frac{\exp(-\frac{1}{2}(\mathbf{x}_i - \mathbf{1}\delta_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \mathbf{1}\delta_i))}{(2\pi)^{n/2} \sqrt{|\boldsymbol{\Sigma}|}}. \end{aligned} \quad (16)$$

Let us define $\bar{x}_i = \sum_{j=1}^n \mathbf{x}_{ij}/n$. The MSE of the maximum likelihood estimator is:

$$\text{MSE}_{\text{MLE}} = \iint (\delta_i - \bar{x}_i)^2 N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma}) p(\delta_i) d\mathbf{x}_i d\delta_i$$

Consider that $(\delta_i - \bar{x}_i)^2 = \left(\delta_i - \frac{1}{n} \mathbf{1}^T \mathbf{x}_i\right)^2$ where $\frac{1}{n} \mathbf{1}^T$ is a linear transformation of the variable \mathbf{x}_i . From the properties of the Normal distribution, it follows that

$$\int \left(\delta_i - \frac{1}{n} \mathbf{1}^T \mathbf{x}_i\right)^2 N(\mathbf{x}_i; \mathbf{1}\delta_i, \boldsymbol{\Sigma}) d\mathbf{x}_i = \frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1}$$

and since

$$\int \left(\frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1}\right) p(\delta_i) d\mathbf{x}_i d\delta_i = \frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1},$$

we derive the first result.

Proof of Proposition 43 The MSE of the shrunk estimator can be obtained in a similar way. First observe that

$$\begin{aligned} &(\delta_i - w\bar{x}_i - (1-w)\delta_0)^2 \\ &= w^2 (\delta_i - \bar{x}_i)^2 + (1-w)^2 (\delta_i - \delta_0)^2 \\ &\quad + 2w(1-w) (\delta_i - \bar{x}_i) (\delta_i - \delta_0) \end{aligned}$$

and its expected value w.r.t. $N(\mathbf{x}_i; \delta_i, \sigma_n^2) p(\delta_i)$ is:

$$\begin{aligned} &\int \left[w^2 \sigma_n^2 + (1-w)^2 (\delta_i - \delta_0)^2 \right] p(\delta_i) d\delta_i \\ &= w^2 \sigma_n^2 + (1-w)^2 \sigma_0^2. \end{aligned} \quad (17)$$

where we have denoted $\sigma_n^2 = \frac{1}{n^2} \mathbf{1}^T \boldsymbol{\Sigma} \mathbf{1}$.

References

- Benavoli, Alessio, Corani, Giorgio, Mangili, Francesca, Zaffalon, Marco, and Ruggeri, Fabrizio. A Bayesian Wilcoxon signed-rank test based on the Dirichlet process. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pp. 1026–1034, 2014.
- Bouckaert, Remco R. Choosing between two learning algorithms based on calibrated tests. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 51–58, 2003.
- Carpenter, Bob, Lee, Daniel, Brubaker, Marcus A, Riddell, Allen, Gelman, Andrew, Goodrich, Ben, Guo, Jiqiang, Hoffman, Matt, Betancourt, Michael, and Li, Peter. Stan: A probabilistic programming language. *Journal of Statistical Software*, in press, 2016.
- Corani, Giorgio and Benavoli, Alessio. A Bayesian approach for comparing cross-validated algorithms on multiple data sets. *Machine Learning*, 100(2):285–304, 2015.
- Demšar, Janez. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, 2006.
- Gelman, Andrew. Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian analysis*, 1(3):515–534, 2006.
- Gelman, Andrew, Carlin, J. B, Stern, H. S, and Rubin, D. B. *Bayesian data analysis*, volume 2. Taylor & Francis, 2014.
- Juárez, Miguel A and Steel, Mark FJ. Model-based clustering of non-Gaussian panel data based on skew-t distributions. *Journal of Business & Economic Statistics*, 28(1):52–66, 2010.
- Kohavi, Ron. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence-Volume 2*, pp. 1137–1143. Morgan Kaufmann Publishers Inc., 1995.
- Kruschke, John. *Doing Bayesian data analysis: A tutorial with R, Jags and Stan*. Academic Press, 2015.
- Kruschke, John K. Bayesian estimation supersedes the t-test. *Journal of Experimental Psychology: General*, 142(2):573, 2013.
- Lacoste, Alexandre, Laviolette, François, and Marchand, Mario. 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)*, pp. 665–675, 2012.
- Lecoutre, Bruno and Poitevineau, Jacques. *The Significance Test Controversy Revisited*. Springer, 2014.
- Murphy, Kevin P. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- Nadeau, Claude and Bengio, Yoshua. Inference for the generalization error. *Machine Learning*, 52(3):239–281, 2003.
- Wagenmakers, Eric-Jan. A practical solution to the pervasive problems of p values. *Psychonomic bulletin & review*, 14(5):779–804, 2007.
- Witten, Ian H, Frank, Eibe, and Hall, Mark. *Data Mining: Practical machine learning tools and techniques (third edition)*. Morgan Kaufmann, 2011.
- Zhang, Harry, Jiang, Liangxiao, and Su, Jiang. Hidden naive Bayes. In *Proceedings of the National Conference on Artificial Intelligence*, number 2, pp. 919. Menlo

Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2005.