
Process-Oriented Evaluation: The Next Step

Pedro Domingos

Artificial Intelligence Group
Instituto Superior Técnico
Lisbon 1049-001, Portugal
pedrod@gia.ist.utl.pt
<http://www.gia.ist.utl.pt/~pedrod>

Abstract

Methods to avoid overfitting fall into two broad categories: data-oriented (using separate data for validation) and representation-oriented (penalizing complexity in the model). Both have limitations that are hard to overcome. We argue that fully adequate model evaluation is only possible if the search process by which models are obtained is also taken into account. To this end, we recently proposed a method for *process-oriented evaluation (POE)*, and successfully applied it to rule induction (Domingos, 1998a). However, for the sake of simplicity this treatment made two rather artificial assumptions. In this paper the assumptions are removed, and a simple formula for model evaluation is obtained. Empirical trials show the new, better-founded form of POE to be as accurate as the previous one, while further reducing theory sizes.

1 INTRODUCTION

Overfitting avoidance is a central problem in machine learning and statistics (Cheeseman & Oldford, 1994). If a learner is sufficiently powerful, whatever representation and search methods it uses, it must guard against selecting a model that fits the training data well but captures the underlying phenomenon poorly. Current methods to address it fall into two broad categories. *Data-oriented evaluation* uses separate data to learn and validate models, and includes methods like cross-validation (Breiman, Friedman, Olshen, & Stone, 1984; Stone, 1974), the bootstrap (Efron & Tibshirani, 1993), and reduced-error pruning (Brunk & Pazzani, 1991). It has several disadvantages: it is often computationally intensive, reduces the data available for learning, can be unreliable if the val-

idation set is small, and is itself prone to overfitting if a large number of models is compared (Ng, 1997). *Representation-oriented evaluation* seeks to avoid these problems by using the same data for training and validation, but *a priori* penalizing some models. Bayesian approaches in general fall into this category (Cheeseman, 1990; MacKay, 1992; Chickering & Heckerman, 1997). Representation-oriented measures typically contain two terms, one reflecting fit to the data, and one penalizing model complexity (Akaike, 1978; Schwarz, 1978; Wallace & Boulton, 1968; Rissanen, 1978; Moody, 1992). This approach is only appropriate when the simpler models are truly the more accurate ones, and there is mounting evidence that this is typically not the case ((Domingos, 1998b, 1997; Schuurmans, Ungar, & Foster, 1997; Lawrence, Giles, & Tsoi, 1997; Webb, 1996; Schaffer, 1993; Murphy & Pazzani, 1994), etc.). Structural risk minimization (Vapnik, 1995; Shawe-Taylor, Bartlett, Williamson, & Anthony, 1996; Scheffer & Joachims, 1998) and PAC learning (Kearns & Vazirani, 1994) are representation-oriented methods that seek to bound the difference between training and generalization error using a function of the model space's (effective) dimension. This typically produces bounds that are overly broad, and requires severely restricting the model space.

We believe the limitations of representation-oriented evaluation stem from ignoring the search process by which candidate models¹ are obtained. A learner with an unlimited model space can avoid overfitting as long as it attempts only a limited number of models (even if it is not possible *a priori* to predict which). Intuitively, the more search has been performed to obtain a model, the higher its expected generalization error for a given training-set error. In a recent paper (Domingos, 1998a) we made this intuition precise and applied the resulting formulas to the CN2 rule learner (Clark & Niblett, 1989), obtaining systematic improvements in

¹By "model" we mean model structure *and* parameter values.

generalization error and theory size. However, for the sake of simplicity the treatment in (Domingos, 1998a) made two rather artificial assumptions: that all error rates are *a priori* equally likely, and that a model's generalization error can be roughly estimated by treating all previously-generated models as having similar generalization errors. In this paper we remove these two assumptions, interpret the result, and successfully apply it to CN2.

2 PROCESS-ORIENTED EVALUATION

Suppose learner L_m consists of drawing m hypotheses at random (independently) from some model space, and returning the one with lowest error on a training sample S . Let $h_{m,i}$ be the i th hypothesis generated by L_m . If $h_{m,i}$'s true error rate is $\epsilon_{m,i}$ and S consists of n independently drawn examples, the number of errors $e_{m,i}$ committed by $h_{m,i}$ on S is a binomially distributed variable with parameters n and $\epsilon_{m,i}$:

$$\begin{aligned} p(e_{m,i}|n, \epsilon_{m,i}) &= b(e_{m,i}|n, \epsilon_{m,i}) \\ &= \binom{n}{e_{m,i}} \epsilon_{m,i}^{e_{m,i}} (1 - \epsilon_{m,i})^{n-e_{m,i}} \end{aligned} \quad (1)$$

Let $B(e_{m,i}|n, \epsilon_{m,i})$ be the probability that the number of errors is greater than $e_{m,i}$:

$$B(e_{m,i}|n, \epsilon_{m,i}) = \sum_{i=e_{m,i}+1}^n b(i|n, \epsilon_{m,i}) \quad (2)$$

Notice that this notation is the opposite of the usual notation for a cumulative distribution function (i.e., $B(e|n, \epsilon) = 1 - \text{Binomial_cdf}(e|n, \epsilon)$). It will be more convenient for what follows.

The probability of L_m returning a hypothesis h_m that misclassifies e_m training examples is the probability that at least one of the m hypotheses $h_{m,i}$ makes e_m errors, and all the others make e_m or more errors. Equivalently, it's the probability that all hypotheses $h_{m,i}$ make more than $e_m - 1$ errors, minus the probability that they all make more than e_m errors:

$$\begin{aligned} p(e_m|n, \vec{\epsilon}_m) &= \prod_{i=1}^m B(e_m - 1|n, \epsilon_{m,i}) \\ &\quad - \prod_{i=1}^m B(e_m|n, \epsilon_{m,i}) \end{aligned} \quad (3)$$

where $\vec{\epsilon}_m = (\epsilon_{m,1}, \dots, \epsilon_{m,i}, \dots, \epsilon_{m,m})$. By Bayes' theorem:

$$p(\vec{\epsilon}_m|n, e_m) \propto p(\vec{\epsilon}_m) p(e_m|n, \vec{\epsilon}_m) \quad (4)$$

Let $h_{m,c}$ be the hypothesis with lowest error (i.e., the "chosen" hypothesis, so that learner L_m returns $h_{m,c}$ and $e_m = e_{m,c}$). Our goal is to predict $h_m = h_{m,c}$'s true error rate $\epsilon_m = \epsilon_{m,c}$ from e_m . For this purpose, we marginalize Equation 4 over all the $e_{m,i}$ save $e_{m,c}$:²

$$p(\epsilon_{m,c}|n, e_m) \propto \int_{\vec{\epsilon}_m \setminus c} p(\vec{\epsilon}_m) p(e_m|n, \vec{\epsilon}_m) d\vec{\epsilon}_m \quad (5)$$

where the integral is multiple, over all components of $\vec{\epsilon}_m$ save $\epsilon_{m,c}$. The expected value of $\epsilon_{m,c}$ can now be computed by integration:

$$E[\epsilon_{m,c}|n, e_m] = \frac{\int_0^1 \epsilon_{m,c} p(\epsilon_{m,c}|n, e_m) d\epsilon_{m,c}}{\int_0^1 p(\epsilon_{m,c}|n, e_m) d\epsilon_{m,c}} \quad (6)$$

Let:

$$f = \int_0^1 p(\epsilon_{m,i}) b(e_m|n, \epsilon_{m,i}) d\epsilon_{m,i} \quad (7)$$

$$F = \int_0^1 p(\epsilon_{m,i}) B(e_m|n, \epsilon_{m,i}) d\epsilon_{m,i} \quad (8)$$

$$E_b[\epsilon_{m,c}] = \frac{\int_0^1 \epsilon_{m,c} p(\epsilon_{m,c}) b(e_m|n, \epsilon_{m,c}) d\epsilon_{m,c}}{\int_0^1 p(\epsilon_{m,c}) b(e_m|n, \epsilon_{m,c}) d\epsilon_{m,c}} \quad (9)$$

$$E_B[\epsilon_{m,c}] = \frac{\int_0^1 \epsilon_{m,c} p(\epsilon_{m,c}) B(e_m|n, \epsilon_{m,c}) d\epsilon_{m,c}}{\int_0^1 p(\epsilon_{m,c}) B(e_m|n, \epsilon_{m,c}) d\epsilon_{m,c}} \quad (10)$$

Substituting Equation 3 into 5 and 5 into 6, using the assumption of independent hypotheses, and assuming the same prior $p(\epsilon_{m,i})$ for all hypotheses, we obtain the following expression:

²This is where we previously assumed that $\forall_i \epsilon_{m,i} = \epsilon_{m,c}$ and dropped the prior $p(\vec{\epsilon}_m)$.

$$E[\epsilon_{m,c}|n, e_m] = \frac{f(F+f)^{m-1}}{(F+f)^m - F^m} E_b[\epsilon_{m,c}] - \frac{F[(F+f)^{m-1} - F^{m-1}]}{(F+f)^m - F^m} E_B[\epsilon_{m,c}] \quad (11)$$

For all but the smallest n , $F \gg f$ (Equations 7, 8, 1 and 2). Thus, using the binomial expansion of $(F+f)^m$ we obtain that $(F+f)^m - F^m \simeq m f F^{m-1}$, $(F+f)^{m-1} - F^{m-1} \simeq (m-1) f F^{m-2}$, and $(F+f)^{m-1} \simeq F^{m-1}$. Substituting these into Equation 11 and simplifying, we obtain:

$$E[\epsilon_{m,c}|n, e_m] = \frac{E_b[\epsilon_{m,c}] + (m-1)E_B[\epsilon_{m,c}]}{m} \quad (12)$$

Let $\epsilon_{m,c}^{ML} = e_m/n$ be the maximum likelihood estimate of $\epsilon_{m,c}$. For sufficiently large n , $E_b[\epsilon_{m,c}] \simeq \epsilon_{m,c}^{ML}$ (Equation 9, given a well-behaved prior $p(\epsilon_{m,c})$, i.e., as long as $p(\epsilon_{m,c}) \neq 0$ in the neighborhood of $\epsilon_{m,c} = e_m/n$). Let $\epsilon_{m,c}^{Prior} = \int_0^1 \epsilon_{m,c} p(\epsilon_{m,c}) d\epsilon_{m,c}$ be the prior expected value of $\epsilon_{m,c}$. Suppose a beta or similarly bell-shaped prior is used (Bernardo & Smith, 1994); this is what makes intuitive sense for error rates. In general e_m/n (the inflection point of $B(e_m|n, \epsilon_{m,c})$ as a function of $\epsilon_{m,c}$) will fall below $\epsilon_{m,c}^{Prior}$ (the peak of the prior), since e_m will tend to zero as more hypotheses are generated and the one with lowest error selected. Then, for sufficiently large n , $B(e_m|n, \epsilon_{m,c}) \simeq 1$ over the entire range where $p(\epsilon_{m,c})$ is significantly greater than zero (leaving out only the left tail of the distribution), and $E_B[\epsilon_{m,c}] \simeq \epsilon_{m,c}^{Prior}$ (Equation 10). Making these substitutions we finally obtain (omitting the c indexes, since $\epsilon_m = \epsilon_{m,c}$):

$$E[\epsilon_m|n, e_m] = \frac{\epsilon_m^{ML} + (m-1)\epsilon_m^{Prior}}{m} \quad (13)$$

This formula is quite similar to the well-known Laplace correction or m -estimate (Cestnik, 1990; Good, 1965). Its role for the number of hypotheses is similar to the m -estimate's role for the number of examples. The m -estimate gradually changes from the maximum likelihood estimate to the prior as the number of examples decreases; similarly, Equation 13 gradually uncovers the prior as the number of hypotheses generated increases. The intuitive meaning of Equation 13 is clear: when a learner generates a series of hypotheses and returns the one with lowest training-set error, the more hypotheses it generates the less sure we are that the observed error corresponds to the true error, and the

more weight should be given to the *a priori* expected error.

This result is intuitively satisfying, because it gives a mathematical basis for increasing model uncertainty as the amount of search performed increases. However, Equation 13 as it stands is of limited practical use, because it converges very rapidly to ϵ_m^{Prior} as more independent hypotheses are generated. As a result, for all but the earliest few hypotheses, the error estimate $E[\epsilon_m|n, e_m]$ is quite insensitive to the empirical error ϵ_m^{ML} . This effect, however, is at least partly due to the fact that hypothesis dependences are being ignored, and as a result the empirical error of one hypothesis carries no information about the true error of another. In particular, only the empirical error of the chosen hypothesis carries information about its true error, resulting in the chosen hypothesis' expected error being the unalloyed prior in all *a priori* possible situations where the minimum empirical error is not the chosen hypothesis' (Equation 3). In practical learners, on the other hand, the hypotheses generated are typically very strongly dependent. Thus, in general, all the empirical errors observed will carry information about the true error of the chosen hypothesis, and Equation 13 should converge correspondingly slower to the prior term ϵ_m^{Prior} . We propose to model this by replacing m in Equation 13 by a slower-growing function of m , which can be thought of as the "effective number of independent hypotheses attempted." For example, attempting ten hypotheses with given dependences between them may be equivalent (with respect to the convergence of Equation 13 to ϵ_m^{Prior}) to attempting two independent hypotheses. Thus, Equation 13 provides a simple way of combining data-oriented, representation-oriented and process-oriented information when estimating generalization error: ϵ_m^{ML} is the data-oriented component (the model's empirical error), ϵ_m^{Prior} is the representation-oriented component (a function of the model's form), and m is the process-oriented component (a function of the search process that led to the model).

3 AN APPLICATION: RULE INDUCTION

Most rule induction systems employ a set covering or "separate and conquer" search strategy (Michalski, 1983; Clark & Niblett, 1989). Rules are induced one at a time, and each rule starts with a training set composed of the examples not covered by any previous rules. A rule is induced by adding conditions one at a time, starting with none (i.e., the rule initially covers the entire instance space). The next condition to add is chosen by attempting all possible conditions.

Conditions on symbolic attributes are typically of the form $a_i = v_{ij}$, where v_{ij} is a possible value of attribute a_i . Conditions on numeric attributes are typically of the form $a_i \leq v_{ij}$ or $a_i > v_{ij}$, where the thresholds v_{ij} are usually values of the attribute that appear in the training set. In the beam search process used by many rule learners, at each step the best b versions of the rule according to some evaluation function are selected for further specialization. AQ (Michalski, 1983) continues adding conditions until the rule is “pure” (i.e., until it covers examples of only one class). This can lead to severe overfitting. The latest version of the CN2 system (Clark & Niblett, 1989; Clark & Boswell, 1991) uses a simple and effective Bayesian method to combat this: induction of a rule stops when no specialization improves its error rate, and the latter is computed using a *Laplace correction* or *m-estimate*. If n_r is the number of examples covered by a rule r , and e_r is the number of those examples it misclassifies, the conventional estimate of the rule’s error rate is e_r/n_r , but its m-estimate is:

$$\hat{\epsilon}_r = \frac{e_r + m\epsilon_0}{n_r + m} \quad (14)$$

where ϵ_0 is the rule’s *a priori* error, which CN2 takes to be the error obtained by random guessing if all classes are equally likely: $\epsilon_0 = (c - 1)/c$, where c is the number of classes. This prior value is given a weight of m examples (i.e., the behavior of Equation 14 is equivalent to having m additional examples covered by the rule, one of each class). CN2 uses $m=c$. As conditions are added, the rule covers fewer and fewer examples, and $\hat{\epsilon}_r$ tends to ϵ_0 . Thus a rule making more misclassifications may be preferred if it covers more examples, causing induction to stop earlier and reducing overfitting. Clark and Boswell (1991) found this version of CN2 to be more accurate than C4.5 (Quinlan, 1993) on 10 of the 12 benchmark datasets they used for testing. However, this scheme ignores that, as more and more conditions are attempted, the probability of finding one that appears to reduce the rule’s error merely by chance increases. This will lead the m-estimate to underestimate the chosen condition’s true error, and CN2 to overfit. The upward correction made to ϵ_r should increase with the number of conditions attempted. The process-oriented evaluation framework described in the previous section allows us to do this in a systematic way, as follows.

Equation 13 can be used to compare the hypotheses returned by k learners $L_1, \dots, L_m, \dots, L_k$, and choose the one with lowest predicted error. It can also be used to compare successive stages of the same learner, by taking L_{m_2} to be the result of continuing the search of learner L_{m_1} ($m_1 < m_2$) with $m_2 - m_1$ more hypothe-

ses. In particular, the successive stages can be the successive versions of a rule returned by CN2 or a similar “separate and conquer” rule learner. A natural choice for the prior expected error $\epsilon_{m,c}^{Prior}$ for all rule versions is the default error rate, obtained by always predicting the most frequent class in the training set. The choice of slower-growing function of m is less obvious. One possibility is $m' = \log m$ (for $m > 1$), based on an analogy with decision tree induction. When learning a tree using an algorithm like CART (Breiman et al., 1984), ID3 (Quinlan, 1986) or C4.5 (Quinlan, 1993), each new hypothesis is obtained by modifying the previous one in only a fraction of the instance space (the fraction corresponding to the node currently being expanded), and this fraction becomes exponentially smaller as induction progresses. Only an entire new level of the decision tree corresponds to an entirely new hypothesis. Since the depth of the tree grows on average with the logarithm of the number of nodes, we can take the equivalent number of independent hypotheses attempted m' to be proportional to the logarithm of the total number of hypotheses attempted m . Since a rule corresponds to a path through a decision tree, both in its content and in the way it is induced by a system like CN2, we can apply a similar line of reasoning to the number of rules attempted.³

Let each hypothesis be one version of the rule attempted during the beam search. Equation 13 does not need to be computed for every rule version generated during the beam search. This would introduce a preference for adding some conditions instead of others, which is unlikely to produce good results unless there is domain knowledge supporting such preferences. Instead, Equation 13 can be computed only once for each round. One round consists of generating every possible one-step specialization of each rule version in the beam, and selecting the b best. Thus, if there are a attributes and v is the maximum number of values of any attribute (in the worst case, $v = n$ for numeric attributes), one round corresponds to $O(bav)$ rule versions. Let m_k be the total number of rule versions generated up to, and including, round k . Round 1 consists of the initial rule with no conditions, and $m_1 = 1$. Induction stops when $E[\epsilon_{m_k} | n_{m_k}, e_{m_k}] \geq E[\epsilon_{m_{k-1}} | n_{m_{k-1}}, e_{m_{k-1}}]$, for $k > 1$.

4 EMPIRICAL STUDY

In order to test the effectiveness of process-oriented evaluation, default and process-oriented versions of

³In the experiments described below, the results were not sensitive to the base of the logarithms used. Base 2, base e and base 10 all yielded practically indistinguishable error rates and theory sizes. The results reported are for base 2.

Table 1: Empirical results: error rates and theory sizes of default CN2 and CN2 with two versions of process-oriented evaluation (CN2-POE1 and CN2-POE2).

Dataset	Error rate			Theory size		
	CN2	CN2-POE1	CN2-POE2	CN2	CN2-POE1	CN2-POE2
Breast	30.0±1.4	29.7±1.4	30.3±1.3	114.5±2.4	58.7±2.6	104.9±2.6
Echocardio	32.7±1.2	32.3±1.3	31.2±1.1	42.9±1.2	35.4±2.1	39.2±1.3
Glass	39.0±1.5	38.3±1.7	39.1±1.4	51.8±1.0	54.7±1.1	45.2±1.0
HeartC	20.8±0.8	22.5±0.8	22.4±0.8	57.8±0.9	52.0±1.0	52.6±1.0
HeartH	22.4±1.1	21.8±1.3	21.9±1.1	69.2±1.5	60.3±1.4	58.9±1.1
Hepatitis	21.2±0.9	19.2±1.3	18.8±1.1	40.2±1.7	34.0±1.3	34.4±1.1
Lympho	21.4±1.1	24.1±1.1	23.4±1.2	39.5±0.7	38.7±1.0	32.8±1.1
Soybean	19.5±1.0	19.4±1.0	22.9±1.2	116.7±2.3	110.9±3.1	97.7±1.7
Thyroid	4.1±0.2	3.8±0.2	4.0±0.2	97.5±2.0	104.8±2.0	83.4±2.6
Tumor	60.1±1.0	65.1±1.3	60.0±1.2	302.8±4.6	273.9±4.4	241.6±3.9
Voting	4.8±0.4	4.3±0.3	4.3±0.3	61.7±2.9	49.6±2.5	33.2±1.7

CN2 were compared on the benchmark datasets previously used by Clark and Boswell (1991).⁴ The process-oriented versions were implemented by adding the necessary facilities to the CN2 source code. Details of the earlier version of POE and its implementation can be found in (Domingos, 1998a). CN2’s Laplace estimates are still used to choose the best b specializations in each round. This is preferable to using uncorrected estimates, since as implemented POE has no preference between hypotheses within the same round, and this is also a factor in avoiding overfitting. However, the Laplace correction distorts the value of ϵ_m^{ML} used in Equation 13. This will be particularly pronounced when there are many classes, since CN2 uses $m = c$. In order to minimize this problem, $m = 2$ was used with POE.⁵

The experimental procedure of (Clark & Boswell, 1991) was followed. Each dataset was randomly divided into 67% for training and 33% for testing, and the error rate and theory size (total number of conditions) were measured for default CN2, CN2-POE1 (the earlier version) and CN2-POE2 (the version described in this paper). This was repeated 20 times. The average results and their standard deviations are shown in Table 1;⁶ the results for CN2 and CN2-POE1 are from (Domingos, 1998a).

Compared to CN2-POE1, CN2-POE2 roughly main-

tains accuracy (lower error in five datasets, higher in five, same in one; 0.2% lower error on average) while reducing theory size in most datasets (lower in seven, higher in four, 4.5 fewer conditions on average). This indicates that Equation 13 is successfully deleting unnecessary conditions that the previous method retained. Being in closed form, Equation 13 is also much more efficient to evaluate than the integrals in (Domingos, 1998a).

These results are obviously very preliminary. A version of POE that takes CN2’s search process into account in more detail is currently being developed. We plan to apply it to the datasets above and study its behavior in more detail, using those datasets and synthetic ones.

5 RELATED WORK

The literature on model selection and error estimation is very large, and we will not attempt to review it here. The incompleteness of representation-oriented evaluation was noted 20 years ago by Pearl (1978):

It would, therefore, be more appropriate to connect credibility with the nature of the selection procedure rather than with properties of the final product. When the former is not explicitly known . . . simplicity merely serves as a rough indicator for the type of processing that took place prior to discovery.

Huber (St. Amant & Cohen, 1997; Huber, 1994) expresses thus the need for process-oriented evaluation:

Data analysis is different from, for example, word processing and batch programming: the correctness of the end product cannot be

⁴With the exception of pole-and-cart, which is not available in the UCI repository (Blake, Keogh, & Merz, 1998).

⁵Simply changing $m = c$ to $m = 2$ in default CN2 does not change its performance on the datasets used.

⁶There are some differences between CN2’s results and those reported in (Clark & Boswell, 1991). This may be due to the fact that the default version of CN2 uses a beam size of 5, whereas Clark and Boswell used $b = 20$. The distribution version of CN2 may also differ from the one used in (Clark & Boswell, 1991).

checked without inspecting the path leading to it.

Several pieces of previous work take into account the number of hypotheses being compared, and so can be considered early steps towards process-oriented evaluation. This includes notably systems that use Bonferroni corrections when testing significance (e.g., (Kass, 1980; Gaines, 1989; Jensen & Schmill, 1997); see also (Miller, 1981; Klockars & Sax, 1986; Westfall & Wolfinger, 1997)). A key difference between these systems and what is proposed here is that they require a somewhat arbitrary choice of significance threshold, while this paper directly attempts to optimize the end goal (expected generalization error). Also, the Bonferroni correction does not take hypothesis dependencies into account, while the present framework offers (at least in principle) a way of doing so.

Quinlan and Cameron-Jones's (1995) "layered search" method for automatically selecting CN2's beam width can also be considered a form of process-oriented evaluation. While layered search and the approach proposed here have similar aims, their biases differ: layered search limits the search's width, while the present method limits its length. The latter may be more effective in reducing the fragmentation and small disjuncts problems (Pagallo & Haussler, 1990; Holte, Acker, & Porter, 1989). The assumptions made here are also clearer than those implicit in Quinlan and Cameron-Jones's (1995) measure.

Freund (1998) recently proposed a form of process-oriented evaluation that is closer to the PAC-learning framework. It is an extension of the statistical query model (Kearns, 1993) that attempts to obtain tighter bounds on generalization error by considering the tree of queries that the learner could make. While the general algorithm to obtain these bounds has exponential computational cost in the number of queries made, Freund proposes a specialized version for algorithms based on local search (e.g., CN2) that is more efficient, at the price of loosened bounds. How tight the bounds will be in either case is still an open question; no empirical testing of Freund's (1998) method has been carried out so far. These bounds could be used for model selection by preferring the model with the lowest upper bound (for given parameters). However, as with Bonferroni corrections, the result will in general depend on the choice of parameters, for which there is no clear criterion. While the approach proposed in the present paper directly obtains an estimate of the generalization error, it would also be useful to have a confidence interval for it, and Freund's (1998) method may be a path to it.

Evaluating models that are the result of a search process, not just of fitting the parameters of a pre-determined structure, has traditionally not been a concern of statisticians. However, this is beginning to change (Chatfield, 1995).

Some of the arguments made here for taking into account the number of hypotheses attempted are made in greater detail in (Jensen & Cohen, 1998) and (Ng, 1997). The present paper goes further in also proposing a method for taking dependences between those hypotheses into account, and in proposing a principled way of combining search process information with more traditional representation-based factors.

6 CONCLUSION

Two main types of model selection are currently available. In *data-oriented evaluation*, a hypothesis's score does not depend on its form or how the hypothesis was found, but only on its performance on the data. In *representation-oriented evaluation*, the score depends on the data and on the hypothesis's form, but not on the search process that led to it. Recently (Domingos, 1998a) we argued that the latter cannot be ignored, and proposed *process-oriented evaluation* (POE). However, in (Domingos, 1998a) we assumed that all models searched had similar true error rates, and that all error rates were equally likely *a priori*. In this paper we removed these assumptions, and derived a simple approximation for the generalization error of the returned hypothesis as a function of the number of hypotheses searched. This approximation is a weighted average of the maximum likelihood estimate of the error and the prior expected error, that increasingly favors the prior as more models are attempted. This approximation gives a mathematical basis to the intuition that model uncertainty should increase with the amount of search conducted.

In the future we plan to: study the statistical properties of Equation 11, in particular when the sample size is not large enough to approximate it by Equation 13; compare the method proposed here with other forms of process-oriented evaluation (e.g., Bonferroni corrections (Jensen & Schmill, 1997) and layered search (Quinlan & Cameron-Jones, 1995)); apply it to other learners; and study methods for accurately estimating the growth of the effective number of hypotheses m' in each of these learners.

References

- Akaike, H. (1978). A Bayesian analysis of the minimum AIC procedure. *Annals of the Institute of Statistical Mathematics*, 30A, 9-14.

- Bernardo, J. M., & Smith, A. F. M. (1994). *Bayesian theory*. New York, NY: Wiley.
- Blake, C., Keogh, E., & Merz, C. J. (1998). *UCI repository of machine learning databases*. Irvine, CA: Department of Information and Computer Science, University of California at Irvine. (<http://www.ics.uci.edu/~mlearn/MLRepository.html>)
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and regression trees*. Belmont, CA: Wadsworth.
- Brunk, C., & Pazzani, M. J. (1991). An investigation of noise-tolerant relational concept learning algorithms. In *Proceedings of the Eighth International Workshop on Machine Learning* (pp. 389-393). Evanston, IL: Morgan Kaufmann.
- Cestnik, B. (1990). Estimating probabilities: A crucial task in machine learning. In *Proceedings of the Ninth European Conference on Artificial Intelligence* (pp. 147-149). Stockholm, Sweden: Pitman.
- Chatfield, C. (1995). Model uncertainty, data mining and statistical inference. *Journal of the Royal Statistical Society A*, 158.
- Cheeseman, P. (1990). On finding the most probable model. In J. Shragar & P. Langley (Eds.), *Computational models of scientific discovery and theory formation* (pp. 73-95). San Mateo, CA: Morgan Kaufmann.
- Cheeseman, P., & Oldford, R. W. (1994). Preface. In P. Cheeseman & R. W. Oldford (Eds.), *Selecting models from data: Artificial intelligence and statistics IV*. New York, NY: Springer-Verlag.
- Chickering, D. M., & Heckerman, D. (1997). Efficient approximations for the marginal likelihood of Bayesian networks with hidden variables. *Machine Learning*, 29, 181-212.
- Clark, P., & Boswell, R. (1991). Rule induction with CN2: Some recent improvements. In *Proceedings of the Sixth European Working Session on Learning* (pp. 151-163). Porto, Portugal: Springer-Verlag.
- Clark, P., & Niblett, T. (1989). The CN2 induction algorithm. *Machine Learning*, 3, 261-283.
- Domingos, P. (1997). Why does bagging work? A Bayesian account and its implications. In *Proceedings of the Third International Conference on Knowledge Discovery and Data Mining* (pp. 155-158). Newport Beach, CA: AAAI Press.
- Domingos, P. (1998a). A process-oriented heuristic for model selection. In *Proceedings of the Fifteenth International Conference on Machine Learning* (pp. 127-135). Madison, WI: Morgan Kaufmann.
- Domingos, P. (1998b). Occam's two razors: The sharp and the blunt. In *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining* (pp. 37-43). New York, NY: AAAI Press.
- Efron, B., & Tibshirani, R. J. (1993). *An introduction to the bootstrap*. New York, NY: Chapman and Hall.
- Freund, Y. (1998). Self bounding learning algorithms. In *Proceedings of the Eleventh Annual Conference on Computational Learning Theory*. Madison, WI: Morgan Kaufmann.
- Gaines, B. R. (1989). An ounce of knowledge is worth a ton of data. In *Proceedings of the Sixth International Workshop on Machine Learning* (pp. 156-159). Ithaca, NY: Morgan Kaufmann.
- Good, I. J. (1965). *The estimation of probabilities: An essay on modern Bayesian methods*. Cambridge, MA: MIT Press.
- Holte, R. C., Acker, L. E., & Porter, B. W. (1989). Concept learning and the problem of small disjuncts. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence* (pp. 813-818). Detroit, MI: Morgan Kaufmann.
- Huber, P. J. (1994). Languages for statistics and data analysis. In P. Dirschedl & R. Ostermann (Eds.), *Computational statistics: Papers collected on the occasion of the Twenty-Fifth Conference on Statistical Computing at Schloss Reisensburg*. Heidelberg: Physica-Verlag.
- Jensen, D., & Cohen, P. R. (1998). Multiple comparisons in induction algorithms. *Machine Learning*. (To appear)
- Jensen, D., & Schmill, M. (1997). Adjusting for multiple comparisons in decision tree pruning. In *Proceedings of the Third International Conference on Knowledge Discovery and Data Mining* (pp. 195-198). Newport Beach, CA: AAAI Press.
- Kass, G. V. (1980). An exploratory technique for investigating large quantities of categorical data. *Applied Statistics*, 29, 119-127.
- Kearns, M. (1993). Efficient noise-tolerant learning from statistical queries. In *Proceedings of the Twenty-Fifth ACM Symposium on the Theory of Computing* (pp. 392-401). New York, NY: ACM Press.

- Kearns, M. J., & Vazirani, U. V. (1994). *An introduction to computational learning theory*. Cambridge, MA: MIT Press.
- Klockars, A. J., & Sax, G. (1986). *Multiple comparisons*. Beverly Hills, CA: Sage.
- Lawrence, S., Giles, C. L., & Tsoi, A. C. (1997). Lessons in neural network training: Overfitting may be harder than expected. In *Proceedings of the Fourteenth National Conference on Artificial Intelligence* (pp. 540-545). Providence, RI: AAAI Press.
- MacKay, D. (1992). Bayesian interpolation. *Neural Computation*, 4, 415-447.
- Michalski, R. S. (1983). A theory and methodology of inductive learning. *Artificial Intelligence*, 20, 111-161.
- Miller, R. G., Jr. (1981). *Simultaneous statistical inference* (2nd ed.). New York, NY: Springer-Verlag.
- Moody, J. E. (1992). The effective number of parameters: An analysis of generalization and regularization in nonlinear learning systems. In J. E. Moody, S. J. Hanson, & R. P. Lippmann (Eds.), *Advances in neural information processing systems 4* (pp. 847-854). San Mateo, CA: Morgan Kaufmann.
- Murphy, P., & Pazzani, M. (1994). Exploring the decision forest: An empirical investigation of Occam's razor in decision tree induction. *Journal of Artificial Intelligence Research*, 1, 257-275.
- Ng, A. Y. (1997). Preventing "overfitting" of cross-validation data. In *Proceedings of the Fourteenth International Conference on Machine Learning* (pp. 245-253). Nashville, TN: Morgan Kaufmann.
- Pagallo, G., & Haussler, D. (1990). Boolean feature discovery in empirical learning. *Machine Learning*, 3, 71-99.
- Pearl, J. (1978). On the connection between the complexity and credibility of inferred models. *International Journal of General Systems*, 4, 255-264.
- Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1, 81-106.
- Quinlan, J. R. (1993). *C4.5: Programs for machine learning*. San Mateo, CA: Morgan Kaufmann.
- Quinlan, J. R., & Cameron-Jones, R. M. (1995). Over-searching and layered search in empirical learning. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence* (pp. 1019-1024). Montréal, Canada: Morgan Kaufmann.
- Rissanen, J. (1978). Modeling by shortest data description. *Automatica*, 14, 465-471.
- Schaffer, C. (1993). Overfitting avoidance as bias. *Machine Learning*, 10, 153-178.
- Scheffer, T., & Joachims, T. (1998). *Estimating the expected error of empirical minimizers for model selection* (Tech. Rep. No. TR 98-9). Berlin, Germany: Computer Science Department, Technical University of Berlin.
- Schuermans, D., Ungar, L. H., & Foster, D. P. (1997). Characterizing the generalization performance of model selection strategies. In *Proceedings of the Fourteenth International Conference on Machine Learning* (pp. 340-348). Nashville, TN: Morgan Kaufmann.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6, 461-464.
- Shawe-Taylor, J., Bartlett, P. L., Williamson, R. C., & Anthony, M. (1996). *Structural risk minimization over data-dependent hierarchies* (Tech. Rep. No. NC-TR-96-053). Egham, UK: Department of Computer Science, Royal Holloway, University of London.
- St. Amant, R., & Cohen, P. R. (1997). Building an EDA assistant: A progress report. In *Preliminary Papers of the Sixth International Workshop on Artificial Intelligence and Statistics* (pp. 501-512). Ft. Lauderdale, FL: Society for Artificial Intelligence and Statistics.
- Stone, M. (1974). Cross-validators choice and assessment of statistical predictions. *Journal of the Royal Statistical Society B*, 36, 111-147.
- Vapnik, V. N. (1995). *The nature of statistical learning theory*. New York, NY: Springer-Verlag.
- Wallace, C. S., & Boulton, D. M. (1968). An information measure for classification. *Computer Journal*, 11, 185-194.
- Webb, G. I. (1996). Further experimental evidence against the utility of Occam's razor. *Journal of Artificial Intelligence Research*, 4, 397-417.
- Westfall, P. H., & Wolfinger, R. D. (1997). Multiple tests with discrete distributions. *American Statistician*, 51, 3-8.