

A Statistical Learning Perspective of Genetic Programming

Nur Merve Amil, Nicolas Bredeche, Christian Gagné,
Sylvain Gelly, Marc Schoenauer, and Olivier Teytaud
Equipe TAO – INRIA Futurs, LRI, Bat. 490,
Université Paris-Sud, 91405 Orsay CEDEX, France
olivier.teytaud@inria.fr

October 26, 2007

Abstract

Code bloat, the excessive increase of code size, is an important issue in Genetic Programming (GP). This paper proposes a theoretical analysis of code bloat in GP from the perspective of statistical learning theory, a well grounded mathematical toolbox for machine learning. By computing the Vapnik-Chervonenkis dimension of the family of programs that can be inferred by a specific setting of GP, it is proved that a parsimonious fitness ensures universal consistency. This means that the empirical error minimization allows convergence to the best possible error when the number of test cases goes to infinity. However, it is also proved that the standard method consisting in putting a hard limit on the program size still results in programs of infinitely increasing size in function of their accuracy. It is also shown that cross-validation or hold-out for choosing the complexity level that optimizes the error rate in generalization also leads to bloat. So a more complicated modification of the fitness is proposed in order to avoid unnecessary bloat while nevertheless preserving universal consistency.

Key words: Statistical learning theory; Code bloat; Symbolic regression.

1 Introduction

This paper is about two important issues in Genetic Programming (GP), that is Universal Consistency (UC) and code bloat. UC consists in the convergence to the optimal error rate with regards to an unknown distribution of examples. A restricted version of UC is consistency, which focus on the convergence to the optimal error rate within a restricted search space. Both UC and consistency are well studied in the field of statistical learning theory. Despite their possible benefits, they have not been widely studied in the field of GP.

Code bloat is the uncontrolled growth of program size that may occur in GP when relying on a variable length representation [13, 14]. This has been identified as a key problem in GP for which there have been several empirical studies. However, very few theoretical studies addressed this issue directly.

The work presented in this paper is intended to provide some theoretical insights on the bloat phenomenon and its link with UC in the context of GP-based learning taking a statistical learning theory perspective [27]. Statistical learning theory provides several theoretical tools to analyze some aspects of learning accuracy. Our main objective consists in performing both an in-depth analysis of bloat as well as providing appropriate solutions to avoid it.

The following section gives an overview of current code bloat theories, describing the results presented in this paper from a GP perspective and providing a short discussion on these results and their benefits for the GP practitioner. Section 2 and 3 formally prove all the aforementioned results about code bloat avoidance and UC and suggest a new approach ensuring both. Then, Section 4 provides some extensions of the previous theoretical results on the use of cross-validation and hold-out methodologies.

Follows some experimental results in Section 5, illustrating the accuracy of the theoretical results. Section 6 finally concludes this paper with a discussion

on the consequences of those theoretical results for GP practitioners and uncover some perspectives of work.

1.1 Code Bloat in Genetic Programming

As quoted from [3, p. 182]:

[...] in 1994, Angeline noted that many of the evolved solutions in Koza's book contained code segments that were extraneous. By extraneous, he meant that if those code segments were removed from the solution, this would not alter the result produced by the solution. Examples of such code would be: (1) $a = a + 0$ or (2) $b = b * 1$.

While bloat is well-defined and can be identified, there are currently no consensual explanations on why it occurs. Indeed, three popular theories can be found in the literature to explain it:

- The *introns* theory states that bloat acts as a protective mechanism in order to avoid the destructive effects of operators once relevant solutions have been found [5, 19, 20]. Introns are pieces of code that have no influence on the fitness: either sub-programs that are never executed, or sub-programs which have no effect;
- The *fitness causes bloat* theory relies on the assumption that there is a greater probability to find a bigger program with the same behavior (i.e. semantically equivalent) than to find a shorter one. Thus, once a good solution is found, programs naturally tend to grow because of fitness pressure [16]. This theory states that code bloat is operator-independent and may happen for any variable length representation-based algorithm. As a consequence, code bloat is not to be limited to population-based stochastic algorithm (such as GP), but may be extended to many algorithms using variable length representation [14];

- The *removal bias* theory states that removing longer sub-programs is more dangerous to do than removing shorter ones (because of possible destructive consequence), so there is a natural bias that benefits to the preservation of longer programs [24].

While it is now considered that each of these theories somewhat capture part of the problem [2], there has not been any definitive global explanation of the bloat phenomenon. At the same time, no definitive practical solution has been proposed that would avoid the drawbacks of bloat (i.e. increasing evaluation time of large trees) while maintaining the good performances of GP on difficult problems. Some common solutions rely either on specific operators (e.g. size-fair crossover [15], or different fair mutation [17]), on some parsimony-based penalization of the fitness [25] or on abrupt limitation of the program size such as the one originally used by Koza [13]. Also, some multi-objective approaches have been proposed ([18, 22, 9, 7, 4]). Some other more particular solutions have been proposed but are not widely used yet [21, 28].

1.2 Structural and Functional Bloat

Although code bloat is not clearly understood, it is yet possible to distinguish at least two kinds of code bloat. We first define *structural bloat* as the code bloat that necessarily takes place when no optimal solution can be approximated by a set of programs with bounded length. In such a situation, optimal solutions of increasing accuracy will also exhibit an increasing complexity (larger programs), as larger and larger code will be generated in order to better approximate the target function. This extreme case of structural bloat has also been demonstrated in [12]. The authors use some polynomial functions of increasing difficulty, and demonstrate that a precise fit can only be obtained through an increased bloat (see also [6] for related issues about problem complexity in GP).

Another form of bloat is the *functional bloat*, which takes place when program

length keeps on growing even though an optimal solution (of known complexity) does lie in the search space. In order to clarify this point, let us use a simple symbolic regression problem defined as follow: given a set \mathcal{S} of test cases, the goal is to find a function f (here, a GP-tree) that minimizes the Mean Square Error (or MSE). If we intend to approximate a polynomial (e.g. $14 * x^2$ with $x \in [0, 1]$), we may observe code bloat since it is possible to find arbitrarily long polynomials that gives the exact solution (e.g. $14x^2 + 0 * x^3 + \dots$), or sequences of polynomials of length growing to ∞ and accuracy converging to the optimal accuracy (e.g. $P_n(x) = 14x^2 + \sum_{i=1}^n \frac{1}{n!i} x^i$). Most of the works cited earlier are in fact concerned with functional bloat which is the most surprising, and the most disappointing kind of bloat.

We will consider various levels of functional bloat: cases where the length of programs found by GP runs to infinity as the number of test cases run to infinity whereas a bounded-length solution exists, and also cases where large programs are found with high probability by GP whereas a small program is optimal.

1.3 Universal Consistency

Another important issue is to study the convergence of the function given by GP, under some sufficient conditions, when the number of test cases goes to infinity, toward the actual function used to generate the test cases. This property is known in statistical learning as *Universal Consistency* (UC). Note that this notion is slightly different from that of universal approximation, to which people usually refer when doing symbolic regression in GP: because polynomial for instance are known to be able to approximate any continuous function, GP search using operators $\{+, *\}$ is also assumed to be able to approximate any continuous function. However, UC is concerned with the behavior of the algorithm when the number of test cases goes to infinity: the existence of a polynomial that

approximates a given function at any arbitrary precision does not imply that any polynomial approximation built from a set of sample points will converge to that given function when the number of points goes to infinity. Precisely, UC can be stated informally as follows (a formal definition will be given later):

Informal definition of symbolic regression from examples: a program is a symbolic regression from examples if it takes as inputs a finite number of examples x_1, \dots, x_n with their labels y_1, \dots, y_n and outputs a program P_n .

Informal definition of universal consistency: a genetic programming tool for symbolic regression from examples is universally consistent if, for all symbolic regression problems, as the number of examples goes to infinity and if examples are independently identically drawn as some random variables (x, y) , the error when applied on unseen examples decreases to the minimal possible one: $Pr(P_n(x) \neq y) \rightarrow 0$.

This definition is not so simple, due to the various dependencies. Consider P_n the program built by GP when n examples are provided. P_n might be stochastic, as it depend on two stochastic elements: (i) the random seed of the GP algorithm (ii) the stochasticity of the generator of examples. The question is not the dependency of P_n on the computation time; for a given genetic programming tool, with its pros and cons and its stopping criterion, we consider only its final output. Whereas the main stream of GP considers the dynamic of the optimization for a fixed n , we here consider the evolution as $n \rightarrow \infty$.

1.4 Results Overview

The UC of GP is investigated in Sections 2 and 3, with a detailed study of structural and functional bloats that might occur when searching a program space with GP. A formal and detailed definition of the space of programs that we consider in this paper is given in Definition 2.1. Note that various definitions could be considered also, the main elements being: i) universal approximation,

and ii) a measure of program-complexity which leads to the finiteness of the VC-dimension for a given program-complexity. Two types of results are derived: i) *UC-related results*, that is whether the probability of misclassification of GP-solutions converges to the optimal value when the number of test cases goes to infinity, and ii) *bloat-related results*, regarding first structural bloat and then functional bloat with various fitness penalization and complexity bounds.

Let us now state precisely, yet informally, the main results of this paper. Section 2 precisely defines the computing machine under examination, and proves the resulting GP search space fulfills the conditions of some standard statistical learning theorems listed in Appendix A. Applying those statistical learning theorems to GP lead to Theorem 2.4, which demonstrates the UC of GP when the fitness measure includes some complexity penalization. Proposition 2.5, a bloat-related theoretical result, unsurprisingly proves that if optimality can not be reached within finite complexity, then converging to the optimal error implies an infinite increase of bloat. And then, Theorem 2.6 proves a negative result about bloat, that is minimizing the MSE alone might lead to bloat even if an optimal function with bounded length belongs to the program search space (i.e. empirical solutions complexity goes to infinity with the sample size).

Then in Section 3, Theorems 2.7 and 3.1 show that it is possible to carefully adjust the parsimony pressure in order to obtain both UC and bounds on the empirical solution complexity (i.e. no bloat). These are the best positive results one could expect considering the previous findings. Section 4 discuss some properties of alternate solutions for complexity penalization, that is cross-validation or hold out, with various pairings of data sets. Note that, though all proofs in Section 2 are stated and proved in the context of binary classification (i.e. find a function from \mathbb{R}^d into $\{0, 1\}$), their generalization to regression (i.e. find a function from \mathbb{R}^d into \mathbb{R}) is straightforward.

2 A First Step Towards No Bloat

The following pages present formal proofs of the results detailed in previous section. Note there are intensive references to elements from statistical learning theory that are presented in the Appendix A. The reader unfamiliar with this theory is advised to read the appendix before reading the present section.

First, Definition 2.1 precisely defines the programs space under examination and Lemma 2.2 carefully shows that it satisfies the hypotheses of Theorems A.1 to A.4 of Appendix A. As stated by Theorem 2.3, this allows the evaluation of the VC-dimension [27] of sets of programs. Then, Theorem 2.4, Proposition 2.5 and Theorem 2.6 are derived from these preliminary results.

It should be noted the mildness of the hypothesis behind Lemma 2.2. We consider any programs of bounded length, working with real variables, provided that the computation time is *a priori* bounded. Usual families of programs in GP verify this hypothesis and much stronger hypothesis. For example, usual tree-based representations avoid loops and therefore all quantities that have to be bounded in lemma below (typically, number of times each operator is used) are bounded for trees of bounded depths. This is also true for direct acyclic graphs. We here deal with a very general case; much better constants can be derived for specific cases, without changing the fundamental results in the sequel of the paper. Lemma 2.2 is necessary because Theorem A.4 deals with the VC-dimension of one program with parameters, whereas we want to deal with families of programs. Lemma 2.2 provides a (tedious, but conceptually simple) generalization of Theorem A.4 to families of programs, by simulating a family of programs by a computing machine that is a general parametric program. Then, Theorem 2.4, Proposition 2.5 and Theorem 2.6 are derived from these preliminary results.

Definition 2.1 (Set of programs studied). *Let $F(n, t, q, m, z)$ be the set of*

functions from \mathbb{R}^{z-m} towards $\{0, 1\}$ which can be computed by a program with a maximum of n lines as follows:

- The program uses at most t operations.
- Each line contains one operation among the followings.

The operations are taken among the following set of instructions:

- Operations $\alpha \mapsto \exp(\alpha)$ (at most q times);
- Operations $+$, $-$, \times , and $/$;
- Jumps conditioned on $>$, \geq , $<$, \leq , and $=$;
- Output 0;
- Output 1;
- Labels for jumps;
- Constants (at most m different);
- Variables (at most z different, with $z \geq m$).

We note $F(n, t, q, m, z)$ as F for short when there is no ambiguity. The parameters n, t, q, m, z are then implicit.

Lemma 2.2 (Computing machine for the set of programs studied). *We note $\log_2(x)$ the integer part (ceil) of $\log(x)/\log(2)$.*

With $F(n, t, q, m, z)$ as defined in the Definition 2.1 above, following the notations of Theorem A.4, there exists a parameterized program h such that $F(n, t, q, m, z)$ is included in $H = \{x \mapsto h(a, x); a \in \mathbb{R}^d\}$, where h is a program with at most:

- $t' = T(n, t, z) = t + t \max(3 + \log_2(n) + \log_2(z), 7 + 3 \log_2(z)) + n(11 + \max(9 \log_2(z), 0) + \max(3 \log_2(z) - 3, 0))$ instructions.

- $q' = q$ exponentials.
- $d' = 1 + m$ dimensions of parameters.

written with the same set of instructions.

Interpretation. This lemma states that a family of programs as defined above is included in the parameterizations of one well-chosen program. Equivalently, we can say that for any $n \in \mathbb{N}$, $t \in \mathbb{N}$, $q \in \mathbb{N}$, $m \in \mathbb{N}$ and $z \in \mathbb{N}$, there exists some program h (constrained by t' , q' and d' as stated above) such that any program in $F(n, t, q, m, z)$ is of the form $x \mapsto h(a, x)$ for some $a \in \mathbb{R}^d$. This replaces a family of programs by one parametric program (i.e. a computing machine), and that will be useful for the evaluation of the VC-dimension of a family of programs by Theorem A.4.

Proof. In order to prove this result, we define below a program h as in theorem A.4 that can emulate any of the programs in $F(n, t, q, m, z)$, with at most $t' = T(n, t, z)$, $q' = q$, and $d' = 1 + m$. Let x be the input variable of dimension $\dim(x) \leq z - m$.

The program goes as follow:

- Label INPUT;
- Initialize $variable(1)$ with value $x(1)$;
- Initialize $variable(2)$ with value $x(2)$;
- ...
- Initialize $variable(\dim(x))$ with value $x(\dim(x))$;
- Label CONSTANTS;
- Initialize $variable(\dim(x) + 1)$ with value a_1 ;

- Initialize $variable(dim(x) + 2)$ with value a_2 ;
- ...
- Initialize $variable(dim(x) + m)$ with value a_m ;
- Label DECODE-INTO-C;
- Operation decode c ;
- Label BEGIN-OF-MAIN-LOOP;
- Label LINE(1);
- Operation $c(1, 1)$ with variables $c(1, 2)$, $c(1, 3)$, and $c(1, 4)$;
- Label LINE(2);
- Operation $c(2, 1)$ with variables $c(2, 2)$, $c(2, 3)$, and $c(2, 4)$;
- ...
- Label LINE(n);
- Operation $c(n, 1)$ with variables $c(n, 2)$, $c(n, 3)$, and $c(n, 4)$;
- Label OUTPUT(0);
- Output 0;
- Label OUTPUT(1);
- Output 1.

We need m real numbers, for parameters, and $4n$ integers $c(.,.)$ for decoding inputs. “Operation decode c ” translates each input real number y (in $[0, 1]$) into a set of four integers $c(.,.)$. This operation can be developed as follows :

1. Let $y \in [0, 1]$;

2. For each $i \in [1, \dots n]$:

- $c(i, 1) = 0$;
- $y = y * 2$;
- If $(y > 1)$ then $c(i, 1) = 1$ and $y = y - 1$;
- $y = y * 2$;
- If $(y > 1)$ then $c(i, 1) = c(i, 1) + 2$ and $y = y - 1$;
- $y = y * 2$;
- If $(y > 1)$ then $c(i, 1) = c(i, 1) + 4$ and $y = y - 1$.

3. For each $j \in [2, 4]$ and $i \in [1, \dots n]$:

- $c(i, j) = 0$;
- $y = y * 2$;
- If $(y > 1)$ then $c(i, j) = 1$ and $y = y - 1$;
- $y = y * 2$;
- If $(y > 1)$ then $c(i, j) = c(i, j) + 2$ and $y = y - 1$;
- $y = y * 2$;
- If $(y > 1)$ then $c(i, j) = c(i, j) + 4$ and $y = y - 1$;
- ...
- $y = y * 2$;
- If $(y > 1)$ then $c(i, j) = c(i, j) + 2^{\log_2(z)-1}$ and $y = y - 1$.

After decoding, $c(i, 1)$ stands for the code instruction, $c(i, 2)$ and $c(i, 3)$ gives the memory addresses where input values can be reached, and $c(i, 4)$ gives the memory address where the output value should be written.

The cost of this is $n(3 + \max(3 \log_2(z), 0))$ “if then”, and $n(3 + \max(3 \log_2(z), 0))$ operators \times , and $n(2 + \max(3(\log_2(z) - 1), 0))$ operators $+$, and $n(3 + \max(3 \log_2(z), 0))$ operators $-$. The overall sum is bounded by $n(11 + \max(9 \log_2(z), 0) + \max(3 \log_2(z) - 3, 0))$.

The result then derives from the rewriting of “operation $c(i, 1)$ with variables $c(i, 2)$, $c(i, 3)$, and $c(i, 4)$ ”. This operation interprets one code instruction and its parameters (as decoded before). It can be developed as follows:

- If $c(i, 1) == 0$ then goto OUTPUT(1);
- If $c(i, 1) == 1$ then goto OUTPUT(0);
- If $c(i, 2) == 1$ then $c = \text{variable}(1)$;
- If $c(i, 2) == 2$ then $c = \text{variable}(2)$;
- ...
- If $c(i, 2) == z$ then $c = \text{variable}(z)$;
- If $c(i, 1) == 7$ then goto LINE(c) (must be encoded by dichotomy with $\log_2(n)$ lines);
- If $c(i, 1) == 6$ then goto EXPONENTIAL(i);
- If $c(i, 3) == 1$ then $b = \text{variable}(1)$;
- If $c(i, 3) == 2$ then $b = \text{variable}(2)$;
- ...
- If $c(i, 3) == z$ then $b = \text{variable}(z)$;
- If $c(i, 1) == 2$ then $a = c + b$;
- If $c(i, 1) == 3$ then $a = c - b$;
- If $c(i, 1) == 4$ then $a = c \times b$;

- If $c(i, 1) == 5$ then $a = c/b$;
- If $c(i, 4) == 1$ then $variable(1) = a$;
- If $c(i, 4) == 2$ then $variable(2) = a$;
- ...
- If $c(i, 4) == z$ then $variable(z) = a$;
- Label END-OF-INSTRUCTION(i).

For each such instruction, at the end of the program, we add the following three lines:

- Label EXPONENTIAL(i);
- $a = \exp(c)$;
- Goto END-OF-INSTRUCTION(i).

Each sequence of the form “if $x=...$ then” (p times) can be encoded by dichotomy with $\log_2(p)$ tests “if ... then goto”. Hence, the expected result. ■

Theorem 2.3 (Finite VC-dimension of the computing machine). *Consider q' , t' and $d' \geq 0$. Let $F(n, t, q, m, z)$ be the set of programs described by Definition 2.1, where $q \leq q'$, $T(n, t, z) \leq t'$, and $1 + m \leq d'$.*

$$\begin{aligned} VCdim(F) &\leq t'^2 d' (d' + 19 \log_2(9d')) \\ &\leq (d'(q' + 1))^2 + 11d'(q' + 1)(t' + \log_2(9d'(q' + 1))) \end{aligned}$$

If $q = 0$ (no exponential) then $VCdim(F) \leq 4d'(t' + 2)$.

Interpretation. The theorem demonstrates that interesting and natural families of programs have finite VC-dimension. Effective methods can associate a VC-dimension to these families of programs.

Proof. Just plug Lemma 2.2 in Theorem A.4. ■

We now consider how to use such results in order to ensure UC. In all of this paper, $\Pr(\cdot)$ denotes probabilities, as the traditional notation $P(\cdot)$ is used for programs.

First, we show why simple empirical risk minimization (i.e. minimizing the error observed without taking into account programs complexity) does not ensure consistency. More precisely, for some distribution of test cases and some i.i.d. (independent identically distributed) sequence of test cases $\{(x_1, y_1), \dots, (x_n, y_n), \dots\}$, there exists P_1, \dots, P_n, \dots such that

$$\forall n \in \mathbb{N}, \forall i \in \{1, 2, \dots, n\} \quad P_n(x_i) = y_i,$$

and however

$$\forall n \in \mathbb{N} \quad \Pr(P_n(x) = y) = 0.$$

This can be proved by considering that x is uniformly distributed in $[0, 1]$ and y is a constant equal to 1. Then, consider P_n , the program that compares its entry to x_1, x_2, \dots, x_n , and outputs 1 if the entry is equal to x_j for some $j \leq n$, and otherwise outputs 0. With probability 1, this program output 0, whereas almost surely the desired output y is 1.

We therefore conclude that minimizing the empirical risk is not enough for ensuring any satisfactory form of consistency. Let's now show that structural risk minimization (i.e. taking into account a penalization for complex structures) can ensure UC and fast convergence when the solution can be written within finite complexity.

Theorem 2.4 (Universal consistency of genetic programming with structural risk minimization). *Consider q_k, t_k, m_k, n_k , and z_k increasing integer sequences. Define \mathcal{F}_k the set of programs with at most t_k lines executed, z_k variables, n_k lines, q_k exponentials, and m_k constants ($\mathcal{F}_k = F(n_k, t_k, q_k, m_k, z_k)$)*

of Definition 2.1) and $\mathcal{F} = \cup_k \mathcal{F}_k$. Then with $q'_k = q_k$, $t'_k = T(n_k, t_k, z_k)$, and $d'_k = 1 + m_k$, define V_k as:

- If $\forall k \ q_k = 0$, then $V_k = 4d'_k(t'_k + 2)$.
- Otherwise, $V_k = (d'_k(q'_k + 1))^2 + 11d'_k(q'_k + 1)(t'_k + \log_2(9d'_k(q'_k + 1)))$.

Now given s test cases, consider $P \in \mathcal{F}$ minimizing $\hat{L}(P) + \sqrt{\frac{32}{s} V(P) \log(es)}$, where $V(P) = V_k$ where k is minimal such that $P \in \mathcal{F}_k$. Then,

- The generalization error, with probability 1, converges to L^* ;
- If one optimal function belongs to \mathcal{F}_k , then for any s and ϵ such that $V_k \log(es) \leq s\epsilon^2/512$, the generalization error with s test cases is larger than $L^* + \epsilon$ with probability at most $\Delta \exp(-s\epsilon^2/128) + 8s^{V_k} \exp(-s\epsilon^2/512)$ where $\Delta = \sum_{j=1}^{\infty} \exp(-V_j)$.

Interpretation. This theorem shows that genetic programming for binary classification, provided that structural risk minimization is performed (i.e. if we optimize an *ad hoc* compromise between complexity of programs and accuracy on empirical data), is universally consistent and verifies some convergence rate properties.

Proof. Just plug Theorem A.5 in Theorem 2.3. Note that Δ is finite because we assumed that the integer sequences were increasing. ■

We now prove the non-surprising fact that if it is possible to approximate the optimal function (the Bayesian classifier) without reaching it exactly, then the complexity of the program runs to infinity as soon as there is convergence of the generalization error to the optimal one.

Proposition 2.5 (Structural bloat in genetic programming). *Consider $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}_3 \subset \dots$, where \mathcal{F}_V is a set of functions from X to $\{0, 1\}$ with VC-*

dimension bounded by V . Consider $(V(s))_{s \in \mathbb{N}}$ a non decreasing sequence of integers and $(P_s)_{s \in \mathbb{N}}$ a sequence of functions such that $P_s \in \mathcal{F}_{V(s)}$.

Define $L_V = \inf_{P \in \mathcal{F}_V} L(P)$ and $V(P) = \inf\{V; P \in \mathcal{F}_V\}$ and suppose that $\forall V \ L_V > L^*$. Then, $(L(P_s) \xrightarrow{s \rightarrow \infty} L^*) \implies (V(P_s) \xrightarrow{s \rightarrow \infty} \infty)$.

Interpretation. This is structural bloat: if the space of programs approximates but does not contain the optimal function and cannot approximate it within bounded size, then bloat occurs. Note that the assumption $\forall V \ L_V > L^*$ holds simultaneously for all V for many distributions, as we consider countable unions of families with finite VC-dimension (e.g. see [8, chap. 18]).

Proof. Define $\epsilon(V) = L_V - L^*$. ϵ is non-increasing as the \mathcal{F}_i are nested. Consider, as assumed in the proposition, that $L(P_s) \xrightarrow{s \rightarrow \infty} L^*$. Consider V_0 a positive integer and let's prove that if s is large enough, then $V(P_s) \geq V_0$. There exists ϵ_0 such that $\epsilon(V_0) > \epsilon_0 > 0$. So, for s large enough, $L(P_s) \leq L^* + \epsilon_0 \implies L_{V_s} \leq L^* + \epsilon_0 \implies L^* + \epsilon(V_s) \leq L^* + \epsilon_0 \implies \epsilon(V_s) \leq \epsilon_0 \implies V_s \geq V_0$. ■

We now show that, even in cases in which an optimal short program exists, the usual procedure (known as the method of Sieves; see also [23]) defined below, consisting in defining a maximum VC-dimension depending upon the sample size (as usually done in practice and as recommended in Theorem A.3) and then using a moderate family of functions, leads to bloat.

Theorem 2.6 (Bloat with the method of Sieves). *Let $\mathcal{F}_1, \dots, \mathcal{F}_k, \dots$ be non-empty sets of functions with finite VC-dimensions V_1, \dots, V_k, \dots , and let $\mathcal{F} = \cup_n \mathcal{F}_n$. Then given s i.i.d. test cases, consider $\hat{P} \in \mathcal{F}_s$ minimizing the empirical risk \hat{L} in \mathcal{F}_s .*

From Theorem A.3, we already know that if $V_s = o(s/\log(s))$ and $V_s \rightarrow \infty$, then

$$\Pr\left(L(\hat{P}) \leq \hat{L}(\hat{P}) + \epsilon(s, V_s, \delta)\right) \geq 1 - \delta$$

and almost surely

$$L(\hat{P}) \rightarrow \inf_{P \in \mathcal{F}} L(P).$$

We now state that if $V_s \rightarrow \infty$, and noting $V(P) = \min\{V_k; P \in \mathcal{F}_k\}$, then $\forall V_0, \delta_0 > 0, \exists \Pr$, a distribution of probability on X and Y , such that $\exists g \in \mathcal{F}_1$ such that $L(g) = L^*$, and for s sufficiently large $\Pr(V(\hat{P}) \leq V_0) \leq \delta_0$.

Interpretation. The result in particular implies that for any V_0 , there is a distribution of test cases such that $\exists g; V(g) = V_1$ and $L(g) = L^*$, with probability 1, $V(\hat{f}) \geq V_0$ infinitely often as s increases. This shows that bloat can occur if we use only an abrupt limit on code size, if this limit depends upon the number of test cases (*a fortiori* if there's no limit). Note that this result, proved thanks to a particular distribution, could indeed be proved for the whole class of classification problems for which the conditional probability of $Y = 1$ (conditionally to X) is equal to $\frac{1}{2}$ in an open subset of the domain.

Proof. The proof is given for the part that is not Theorem A.3. Figure 1 gives an illustration of the proof. Consider $V_0 > 0, \delta_0 > 0, \alpha$ such that $(e\alpha/2^\alpha)^{V_0} \leq \delta_0/2$, s such that $V_s \geq \alpha V_0$, and $d = \alpha V_0$. Consider that x_1, \dots, x_d are d points shattered by \mathcal{F}_d ; such a family of d points exist, by definition of \mathcal{F}_d . Define the probability measure \Pr by the fact that X and Y are independent and $\Pr(Y = 1) = \frac{1}{2}$ and $\Pr(X = x_i) = \frac{1}{d}$. Then, the following holds, with Q the empirical distribution (the average of Dirac masses on the x_i 's):

1. No empty x_i : $\Pr(E_1) \rightarrow 0$, where E_1 is the fact that $\exists i; Q(X = x_i) = 0$, as $s \rightarrow \infty$;
2. No equality: $\Pr(E_2) \rightarrow 0$, where E_2 is the fact that E_1 occurs or $\exists i; Q(Y = 1|X = x_i) = \frac{1}{2}$;
3. The best function is not in \mathcal{F}_{V_0} : $\Pr(E_3|\neg E_2) \leq S(d, d/\alpha)/2^d$, where E_3 is the fact that $\exists g \in \mathcal{F}_{d/\alpha=V_0}; \hat{L}(g) = \inf_{\mathcal{F}_d} \hat{L}$, with $S(d, d/\alpha)$ the

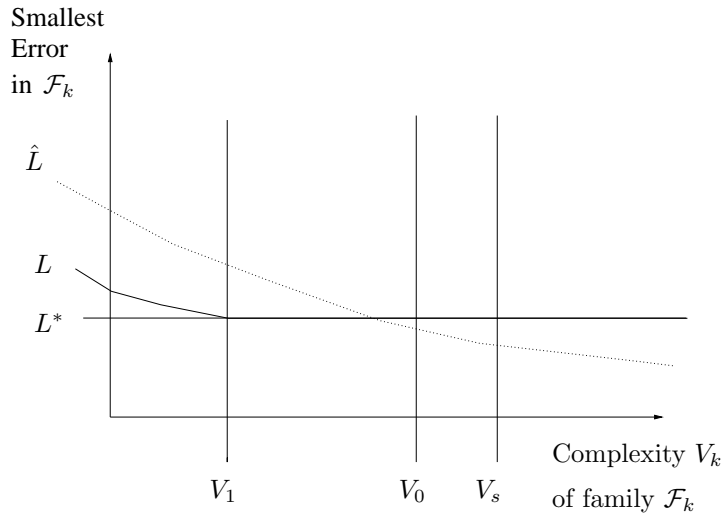


Figure 1: Illustration of the proof of Theorem 2.6 . For V_0 , δ_0 and \Pr fixed, there exists V_1 such that $\inf_{f \in \mathcal{F}_1} L(f) = L^*$, that is to say there exists a simple solution (in \mathcal{F}_1) of error L^* . However, even if the limit V_s ensures that $L(\operatorname{argmin}_{f \in \mathcal{F}_k, V_k \leq V_s} \hat{L})$ is close to L^* , the complexity of the solution can be arbitrarily high ($\geq V_0$ for any V_0)

relevant shattering coefficient, that is the cardinal of $\mathcal{F}_{d/\alpha}$ restricted to $\{x_1, \dots, x_d\}$.

It is now sufficient to rely on classical results. It is well known in the VC-theory that $S(a, b) \leq (ea/b)^b$ (see for example [8, chap.13]), hence $S(d, d/\alpha) \leq (ed/(d/\alpha))^{d/\alpha}$ and $\Pr(E_3 | \neg E_2) \leq (e\alpha)^{d/\alpha} / 2^d \leq \delta_0/2$. If s is sufficiently large to ensure that $\Pr(E_2) \leq \delta_0/2$ (we have proved above that $\Pr(E_2) \rightarrow 0$ as $s \rightarrow \infty$) then

$$\begin{aligned} \Pr(E_3) &\leq \Pr(E_3 | \neg E_2) \Pr(\neg E_2) + \Pr(E_2) \\ &\leq \Pr(E_3 | \neg E_2) + \Pr(E_2) \leq \delta_0/2 + \delta_0/2 \leq \delta_0 \end{aligned}$$

This concludes the proof. ■

Let's now show that it is possible to optimize a compromise between optimality and complexity in an explicit manner (e.g. by replacing 1% precision

with 10 lines of programs or 10 minutes of CPU).

Theorem 2.7 (Bloat-control for regularized empirical risk minimization with relevant VC-dimension). *Let $\mathcal{F}_1, \dots, \mathcal{F}_k, \dots$ be non-empty sets of functions with finite VC-dimensions V_1, \dots, V_k, \dots . Let $\mathcal{F} = \cup_n \mathcal{F}_n$. Consider W a user-defined complexity penalization term. Then, given s test cases, consider $P \in \mathcal{F}_s$ minimizing the regularized empirical risk $\hat{L}(P) = \hat{L}(P) + W(P)$ among \mathcal{F}_s . If $V_s = o(s/\log(s))$ and $V_s \rightarrow \infty$, then $\tilde{L}(\hat{P}) \rightarrow \inf_{P \in \mathcal{F}} \tilde{L}(P)$ almost surely, where $\tilde{L}(P) = L(P) + W(P)$.*

Interpretation. This theorem shows that, using a relevant *a priori* bound on the complexity of the program and adding a user-defined complexity penalization to the fitness, can lead to convergence toward a user-defined compromise [28, 29] between classification rate and program complexity (i.e. we ensure almost sure convergence to a compromise of the form “ λ_1 CPU time + λ_2 misclassification rate + λ_3 number of lines”, where the λ_i are user-defined).

Remark. The drawback of this approach is that we have lost UC and consistency. This means that the misclassification rate in generalization does not converge to the Bayes error in the general case, and in spite of the fact that an optimal program exists, there is not necessarily convergence to its efficiency.

Proof. See Figure 2 for an illustration of the proof. By Theorem A.2,

$$\sup_{P \in \mathcal{F}_s} \left| \hat{L}(P) - \tilde{L}(P) \right| \leq \sup_{P \in \mathcal{F}_s} \left| \hat{L}(P) - L(P) \right| \leq \epsilon(s, V_s)$$

and $\epsilon(s, V_s) \rightarrow 0$ almost surely.

Hence the expected result. ■

Previous results have shown that: i) UC can be reached thanks to usual results of learning theory applied to GP (method of Sieves or structural risk

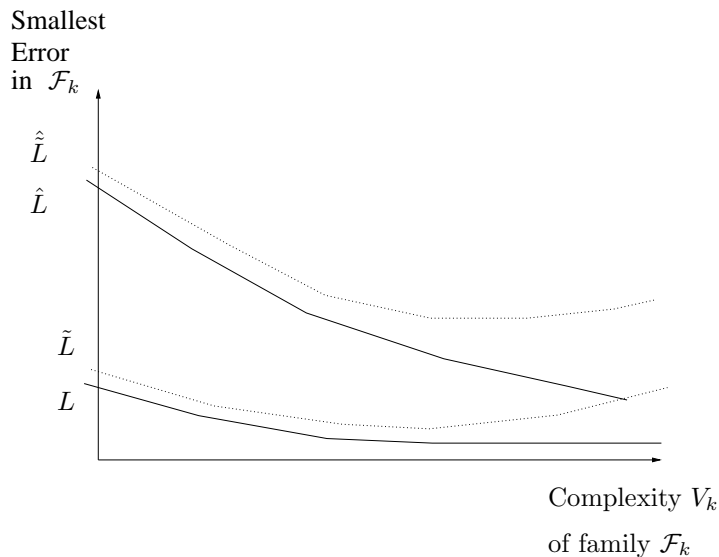


Figure 2: Illustration of the proof of Theorem 2.7. In non-bold plots: smallest \hat{L} (upper plot) or L (lower plot) in \mathcal{F}_k . Bold plots: smallest $\hat{\hat{L}}$ (upper plot) or \tilde{L} (lower plot) in \mathcal{F}_k , i.e. error+penalization. With a larger k , \mathcal{F}_k has a smaller best error \hat{L} , but the penalization is stronger than the gain of error by increasing the complexity

minimization) but sometimes leads to bloat, and ii) bloat can be simply avoided by a strong penalization of size, but this leads to a loss in terms of UC.

In the next section, results will make it possible to present a new approach that combines an *a priori* limit on VC-dimension (i.e. size limit) and a complexity penalization (i.e. parsimony pressure). Indeed, Theorem 3.1 will state that this leads to both UC and convergence to an optimal complexity of the program (i.e. no bloat).

3 Universal Consistency without Bloat

In this section, we consider a more complicated case where the goal is to ensure UC, while simultaneously avoiding non-necessary bloat. This means that an optimal program does exist in a given family of functions and convergence towards the minimal error rate is performed without increasing the program

complexity. This is achieved by: i) merging regularization and bounding of the VC-dimension, and ii) penalization of the complexity (i.e. length) of programs by a penalty term $R(s, P) = R(s)R'(P)$ depending upon the sample size and the program. $R(., .)$ is user-defined and the algorithm looks for a classifier with a small value of both R' and L . In the following, we study both the UC of this algorithm (i.e. $L \rightarrow L^*$) and the no-bloat theorem (i.e. $R' \rightarrow R'(P^*)$ when P^* exists). Note that the bound $V_s = o(\log(s))$ is much stronger than the usual limit used in the method of Sieves (see Theorem 2.6).

Theorem 3.1 (No-bloat theorem). *Let $\mathcal{F}_1, \dots, \mathcal{F}_k, \dots$ with finite VC-dimensions V_1, \dots, V_k, \dots . Let $\mathcal{F} = \cup_n \mathcal{F}_n$. Define $V(P) = V_k$ with $k = \inf\{t | P \in \mathcal{F}_t\}$. Define $L_V = \inf_{P \in \mathcal{F}_V} L(P)$. Consider $V_s = o(\log(s))$ and $V_s \rightarrow \infty$. Consider also that \hat{P}_s minimizes $\hat{L}(P) = \hat{L}(P) + R(s, P)$ in \mathcal{F}_s , and assume that $R(s, .) \geq 0$. Then consistency is attained as whenever $\sup_{P \in \mathcal{F}_{V_s}} R(s, P) = o(1)$, we have that $L(\hat{P}_s) \rightarrow \inf_{P \in \mathcal{F}} L(P)$ almost surely. Note that for well chosen family of functions, $\inf_{P \in \mathcal{F}} L(P) = L^*$. Moreover, assume that $\exists P^* \in \mathcal{F}_{V^*}$ $L(P^*) = L^*$. With $R(s, P) = R(s)R'(P)$ and with $R'(s) = \sup_{P \in \mathcal{F}_{V_s}} R'(P)$, we get the following results:*

1. **Non-asymptotic no-bloat theorem:** *For any $\delta \in]0, 1]$, $R'(\hat{P}_s) \leq R'(P^*) + (1/R(s))2\epsilon(s, V_s, \delta)$ with probability at least $1 - \delta$. This result is in particular interesting for $\epsilon(s, V_s, \delta)/R(s) \rightarrow 0$, which is possible for usual regularization terms as in Theorem A.5 of the Appendix;*
2. **Almost-sure no-bloat theorem:** *If for some $\alpha > 0$, $R(s)s^{(1-\alpha)/2} = O(1)$, then almost surely $R'(\hat{P}_s) \rightarrow R'(P^*)$ and if $R'(P)$ has discrete values (such as the number of instructions in P or many complexity measures for programs) then for s sufficiently large, $R'(\hat{P}_s) = R'(P^*)$;*

3. **Convergence rate:** For any $\delta \in]0, 1]$, With probability at least $1 - \delta$,

$$L(\hat{P}_s) \leq \inf_{P \in \mathcal{F}_{V_s}} L(P) + \underbrace{R(s)R'(s)}_{=o(1) \text{ by hypothesis}} + 2\epsilon(s, V_s, \delta),$$

where

$$\epsilon(s, V, \delta) = \sqrt{\frac{4 - \log(\delta / (4s^{2V}))}{2s - 4}}$$

is such that with probability at least $1 - \delta$, $\epsilon(s, V, \delta) \geq \epsilon(s, V)$ where

$$\epsilon(s, V) = \sup_{f \in \mathcal{F}_V} |\hat{L}(f) - L(f)|,$$

given by Theorem A.1.

Interpretation. Combining a code limitation and a penalization leads to UC without bloat.

Remark. The usual $R(s, P)$ as used in Theorem A.5 or Theorem 2.4 provides consistency and non-asymptotic no-bloat, when this penalization term is used in conjunction with a limitation depending on the sample size s ($P \in \mathcal{F}_s$ with s much more restrictive than in the method of Sieves). A stronger regularization leads to the same results, plus almost sure no-bloat. The asymptotic convergence rate depends upon the regularization. The result is not limited to GP and could be used in other areas. The strongest limitation to this results is not the GP-framework, but the fact that, as shown in Proposition 2.5, the no-bloat results require the fact that $\exists V^* \exists P^* \in \mathcal{F}_{V^*} L(P^*) = L^*$.

Proof. Let's define $\epsilon(s, V) = \sup_{f \in \mathcal{F}_V} |\hat{L}(f) - L(f)|$. For any P , $\hat{L}(\hat{P}_s) + R(s, \hat{P}_s) \leq \hat{L}(P) + R(s, P)$. On the other hand, $L(\hat{P}_s) \leq \hat{L}(\hat{P}_s) + \epsilon(s, V_s)$. So consistency is proved by the following:

$$\begin{aligned} L(\hat{P}_s) &\leq \left(\inf_{P \in \mathcal{F}_{V_s}} (\hat{L}(P) + R(s, P)) \right) - R(s, \hat{P}_s) + \epsilon(s, V_s), \\ &\leq \left(\inf_{P \in \mathcal{F}_{V_s}} (L(P) + \epsilon(s, V_s) + R(s, P)) \right) - R(s, \hat{P}_s) + \epsilon(s, V_s), \\ &\leq \left(\inf_{P \in \mathcal{F}_{V_s}} (L(P) + R(s, P)) \right) + 2\epsilon(s, V_s). \end{aligned}$$

As $\epsilon(s, V_s) \rightarrow 0$ almost surely (see Theorem A.2) and $\inf_{P \in \mathcal{F}_{V_s}} (L(P) + R(s, P)) \rightarrow \inf_{P \in \mathcal{F}} L(P)$, we conclude that $L(\hat{P}_s) \rightarrow \inf_{P \in \mathcal{F}} L(P)$ almost surely.

Let's now focus on the proof of the no-bloat result. By definition of the algorithm, for s sufficiently large to ensure $P^* \in \mathcal{F}_{V_s}$, $\hat{L}(\hat{P}_s) + R(s, \hat{P}_s) \leq \hat{L}(P^*) + R(s, P^*)$, hence with probability at least $1 - \delta$,

$$\begin{aligned} R'(\hat{P}_s) &\leq R'(P^*) + (1/R(s))(L^* + \epsilon(s, V_s, \delta) - L(\hat{P}_s) + \epsilon(s, V_s, \delta)), \\ &\leq R'(V^*) + (1/R(s))(L^* - L(\hat{P}_s) + 2\epsilon(s, V_s, \delta)). \end{aligned}$$

As $L^* \leq L(\hat{P}_s)$, this leads to the non-asymptotic no-bloat version of the theorem.

The almost-sure no-bloat theorem is derived as follows.

$$\begin{aligned} R'(\hat{P}_s) &\leq R'(P^*) + 1/R(s)(L^* + \epsilon(s, V_s) - L(\hat{P}_s) + \epsilon(s, V_s)), \\ &\leq R'(P^*) + 1/R(s)(L^* - L(\hat{P}_s) + 2\epsilon(s, V_s)), \\ &\leq R'(P^*) + 1/R(s)2\epsilon(s, V_s). \end{aligned}$$

All we need is the fact that $\epsilon(s, V_s)/R(s) \rightarrow 0$ almost surely.

For any $\epsilon > 0$, we consider the probability of $\epsilon(s, V_s)/R(s) > \epsilon$, and we sum over $s > 0$. By the Borel-Cantelli lemma¹, the finiteness of this sum is sufficient for the almost sure convergence to 0. The probability of $\epsilon(s, V_s)/R(s) > \epsilon$ is the probability of $\epsilon(s, V_s) > \epsilon R(s)$. By Theorem A.1, this is bounded above by $O(\exp(2V_s \log(s) - 2s\epsilon^2 R(s)^2))$. This has finite sum for $R(s) = \Omega(s^{-(1-\alpha)/2})$.

Let us now consider the convergence rate. Consider s sufficiently large to ensure $L_{V_s} = L^*$. As shown above during the proof of the consistency,

$$\begin{aligned} L(\hat{P}_s) &\leq \inf_{P \in \mathcal{F}_{V_s}} (L(P) + R(s, P)) + 2\epsilon(s, V_s), \\ &\leq \inf_{P \in \mathcal{F}_{V_s}} (L(P) + R(s)R'(P)) + 2\epsilon(s, V_s), \\ &\leq \inf_{P \in \mathcal{F}_{V_s}} (L(P)) + R(s)R'(s) + 2\epsilon(s, V_s). \end{aligned}$$

¹If $\sum_n \Pr(X_n > \epsilon)$ is finite for any $\epsilon > 0$ and $X_n > 0$, then $X_n \rightarrow 0$ almost surely.

So with probability at least $1 - \delta$,

$$L(\hat{P}_s) \leq \inf_{P \in \mathcal{F}_{V_s}} L(P) + R(s)R'(s) + 2\epsilon(s, V_s, \delta).$$

■

4 Extensions

When one tries to learn a relation between x and y , the “true” cost function (typically the mean squared error of a given approximate relation, which is an expectation under some usually unknown law of probability) is generally not available. It is usually replaced by its empirical mean on a finite sample. Minimizing this empirical mean is natural, but this can be done over various families of functions (e.g. trees with depth 1, 2, 3, and so on). Choosing between these various levels is hard. Typically, the empirical mean decreases as the complexity is increased (Figure 3), but this decrease is not generally a decrease of the generalization error, as trees of important depth have usually a very bad generalization error due to overfitting. Therefore, the problem is somewhat multi-objective: there is a conflict between the empirical error and the complexity level. Some GP frameworks including this multi-objective idea are [9, 7, 4].

We have studied above:

- The method consisting in minimizing the empirical error on a given family of functions, that is the error observed on the test cases (leading to bloat (this is an *a fortiori* consequence of Theorem 2.6) without universal consistency (see remark before Theorem 2.4);
- The method consisting in minimizing the empirical error, that is the error observed on the test cases, with a (well chosen) hard bound on the

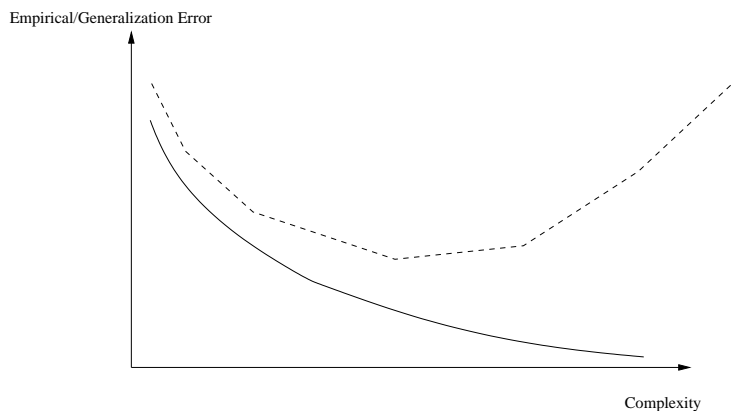


Figure 3: Structural risk minimization or cross-validation to find the ideal complexity level. In the figure, the X-axis denotes the complexity, the Y-axis denotes the empirical error, and the dashed line represents the (usually unknown) generalization error. Structural risk minimization is an ad hoc non-linear compromise between the complexity (X-axis) and the empirical error, the minimization of which ensuring UC (but possibly with bloat). Cross-validation is based on a cross-validation estimate of the dashed line; as shown in this paper, cross-validation also leads to bloat

complexity (method called “method of sieves” [11]; leading to universal consistency but bloat, see Theorem 2.6);

- The method, inspired from structural risk minimization (but slightly adapted against bloat), consisting in minimizing a compromise between the empirical error and the complexity *under the constraint of* a complexity bound including size and computation time (see Theorem 3.1).

All this methods provide a function chosen in a complexity level F_{k^*} . Only the last one provides a estimate of k^* that converges to the smallest value of k^* that allows the optimal generalization error. We now study the following other cases:

- The case in which the level of complexity is chosen through resampling, that is by using cross-validation or hold out (e.g. the program uses cross-validation in order to choose a function on a two-objectives Pareto-set

(empirical error rate vs complexity level));

- The case in which the complexity penalization does not include any time bound but only size bounds.

We mainly conclude that: i) penalization is necessary to avoid bloat, ii) penalization cannot be replaced by cross-validation, and iii) penalization cannot be replaced by hold-out, and must include some time-penalization.

4.1 Validation Sets for Bloat Avoidance

Note UC for Universal Consistency and ERM for Empirical Risk Minimization.

We considered above different cases:

- GP with only “ERM” fitness;
- GP with ERM+bound, typically minimization of the ERM among programs with complexity bounded by some *ad hoc* function of the number of examples (leading to UC + bloat); ERM+bound is a standard tool sometimes called “method of Sieves” [11];
- GP with ERM+penalization+bound, similar to the previous item, but with ERM penalized by an *ad hoc* complexity term (leading to UC without bloat).

One can now consider some other cases:

- Hold out in order to choose between different complexity classes (i.e. on the Pareto-front corresponding to different trade-offs between ERM and complexity, choose the function by hold out);
- Idem through cross-validation.

This section is devoted to these cases. First, let’s consider hold-out for choosing the complexity level. Consider $X_0, \dots, X_N, Y_0, \dots, Y_N$, $2(N + 1)$ samples

(each of them consisting in n examples, i.e. $X_i = (X_{i,1}, X_{i,2}, \dots, X_{i,n})$ and $Y_i = (Y_{i,1}, Y_{i,2}, \dots, Y_{i,n})$), the X_i 's being learning sets, the Y_i 's being (hold out) test sets.

Consider that the function can be chosen in many complexity levels, $F_0 \subset F_1 \subset F_2 \subset F_3 \subset \dots$, where F_0 non-empty and $F_i \neq F_{i+1}$. Note $\hat{L}_k(f)$ the error rate of the function f in the set X_k of examples:

$$\hat{L}_k(f) = \frac{1}{n} \sum_{i=1}^n l(f, X_{k,i})$$

where $l(f, x) = 1$ if f fails on x and 0 otherwise. Define $f_k = \arg \min_{F_k} \hat{L}_k(\cdot)$. In hold-out, after the complete learning, the resulting classifier is f_{k^*} where $k^* = \arg \min_k l_k$ where $l_k = \frac{1}{n} \sum_{i=1}^n l(f_k, Y_{k,i})$.

In the sequel, we assume that $f \in F_k \Rightarrow 1 - f \in F_k$ and that $VCdim(F_k) \rightarrow \infty$ as $k \rightarrow \infty$. The case with hold out leads to different cases, namely:

Greedy case: all X_k 's and Y_k 's are independent; this means that we test separately each complexity level F_k with different learning sets X_k and test sets Y_k .

Case with pairing: X_0 is independent of Y_0 , $\forall k, X_k = X_0$ and $\forall k, Y_k = Y_0$; this means that we use the same learning set for all complexity levels and the same test set for all complexity levels. This case is far more usual.

4.1.1 Case of Greedy Hold-Out

We here consider the choice among complexity levels $0, 1, 2, \dots, N$, though greedy hold-out. For convenience, we define again precisely the framework:

- for $k \in [[0, N(n)]]$, f_k is a function minimizing the error on X_k among F_k ;
- for $k \in [[0, N(n)]]$, l_k is the error of function f_k on Y_k .

k^* is any element of $[[0, N(n)]]$ such that $\forall k \in [[0, N(n)]]$, $l_{k^*} \leq l_k$ (f_k and k^* are not uniquely determined; results below will cover any particular choice). k^*

is the complexity level chosen by HO; bloat occurs when k^* is big whereas small values are enough to reach an optimal generalization error is possible.

$N = N(n)$ is a function of n . We will in the sequel distinguish two cases. First, if $N(n)$ is a constant, then:

- consider the case of an output y independent of the input x , and $\Pr(y = 1) = \Pr(y = 0) = \frac{1}{2}$;
- in that case, an optimal function lies in F_0 - we will show in the following points that, however, k^* will be in general much larger than 0, i.e. bloat occurs;
- then, for any $k \in [[0, N]]$, $n \times l_k$ is a binomial variable with parameters n and $\frac{1}{2}$;
- then, for any $k \in [[1, N]]$ and any $i < k$, $\Pr(N \times l_k = N \times l_i)$ occurs with probability at most $1/\sqrt{n}$ (as the binomial verifies that for some fixed $K > 0$, for any k , $\Pr(B(n, \frac{1}{2}) = k) < K/\sqrt{n}$, see e.g. [8, appendix A.8, lemma A.3]);
- therefore, the probability of the event $E = \{\exists 0 \leq i \neq j \leq N; l_i = l_j\}$ is $O(N^2/n)$;
- therefore, the probability of the event $E = \{\exists 0 \leq i \neq j \leq N; l_i = l_j\}$ goes to 0 as $n \rightarrow \infty$;
- conditionally to the fact that E does not occur, k^* is uniformly distributed on $[[0, N]]$;
- combining the two previous points, $\forall k \in [[0, N]], P(k^* = k) \rightarrow 1/(N + 1)$ as $n \rightarrow \infty$.

This (easier) case of N constant being handled, and leading to the first part of theorem 4.1 below, we now focus on $N(n) \rightarrow \infty$ (which will be the second part of theorem 4.1).

Consider also the case of an output y independent of the input x , and $\Pr(y = 1) = \Pr(y = 0) = \frac{1}{2}$.

- Then, for any k , $n \times l_k$ is binomially distributed (parameters n and $\frac{1}{2}$), as in the previous case.
- For any $i \neq j$, there exists a universal constant K such that $P(l_j < l_i) \geq \frac{1}{2} - K/\sqrt{n}$ (as the binomial verifies that for some fixed $K > 0$, for any k , $\Pr(B(n, \frac{1}{2}) = k) < K/\sqrt{n}$, see e.g. [8, appendix A.8, lemma A.3]).
- For any $C > 0$,

$$P(\exists j \geq C, \forall i < C, l_j < l_i) \geq 1 - \left(1 - \left(\frac{1}{2} - K/\sqrt{n}\right)^C\right)^{N-C}$$

$$\leq 1 - A^{N-C} \text{ for some fixed } A < 1, \text{ and for } n \text{ sufficiently large}$$

- Note $k^*(N)$ the value of k^* chosen for $n = \min\{i \in \mathbb{N}; N(i) = N\}$. The point above shows that, for any $C > 0$, $\sum_{N \in \mathbb{N}} \Pr(k^*(N) < C)$ is finite (as $\Pr(k^* < C)$ decreases exponentially as a function of N for n sufficiently large);
- Therefore, almost surely, thanks to the Borel-Cantelli lemma, $k^*(N) \rightarrow \infty$ as $N \rightarrow \infty$.

We have shown this with *one* distribution, which does not depend upon the number of examples. Moreover, this happens whereas an optimal function lies in F_0 . The result is therefore as follows:

Theorem 4.1 (No bloat avoidance with greedy hold out). *Consider greedy hold-out for choosing between complexity levels $0, 1, \dots, N(n)$.*

- If $N(n)$ is a constant, then for some distribution of examples

$$\forall k \in [[0, N]], P(k^* = k) \rightarrow 1/(N + 1)$$

- If $N(n) \rightarrow \infty$ as $n \rightarrow \infty$, then for some distribution of examples such that an optimal function lies in F_0 , greedy hold-out leads to $k^*(N) \rightarrow \infty$ as $N \rightarrow \infty$ and therefore $\limsup_{n \rightarrow \infty} k^*(n) = \infty$.

4.1.2 Case of Hold-Out with Pairing

We now set $X = X_0 = X_k$ and $Y = Y_0 = Y_k$ for any $k \in [[1, N]]$, i.e. we now consider one learning set and one test set. This is less data-consuming (only up to a $N(n)$ factor) than the greedy case above and could be statistically more significant, thanks to pairing. We below compare it to the previous result.

- Consider $V \in \mathbb{N}$ and $v \in \mathbb{N}$ such that $V = VCdim(F_v)$ with v minimal realizing this condition.
- Consider $A = \{a_1, \dots, a_V\}$, a set of points shattered by F_v .
- Consider a distribution of examples with x uniform on A and y independent of x with $\Pr(y = 1) = \Pr(y = 0) = \frac{1}{2}$.
- Consider \hat{P}_X the empirical distribution associated to the sample X of n examples and \hat{P}_Y the empirical distribution associated to the sample Y of n examples.
- Then, with $E_X = \{\exists i; \hat{P}_X(x = a_i) = 0 \text{ or } \hat{P}_X(y = 1|x = a_i) = \frac{1}{2}\}$, $\lim_n \Pr(E_X) = 0$.
- Then, with $E_Y = \{\exists i; \hat{P}_Y(x = a_i) = 0 \text{ or } \hat{P}_Y(y = 1|x = a_i) = \frac{1}{2}\}$, $\lim_n \Pr(E_Y) = 0$.
- There is at least one function f on A which does not belong to F_{v-1} .

- With probability at least $(1 - \Pr(E_Y))/2^V$, this function f is optimal for $L(\cdot, Y_0)$ (and all optimal functions have the same restriction to A).
- If $k \geq v$, then with probability at least $(1 - \Pr(E_X))/2^V$, $f_k = f$.
- Combining the two probabilities above, as the events are independent, we see that they occur simultaneously with probability at least $p(v, n) = ((1 - \epsilon(v, n))/2^V)^2$, and then $k^* \geq v$, where $\epsilon(v, n) = \max(\Pr(E_X), \Pr(E_Y)) \rightarrow 0$ as $n \rightarrow \infty$. Then:

$$\Pr(k^* \geq v) > p(v, n)$$

- This implies the first result of this section: for any v , $\Pr(k^* \geq v)$ does not go to 0, whereas a function in F_0 is optimal.

This leads to the following proposition.

Proposition 4.2 (Bloat cannot be controlled by hold out with pairing, first result). *For arbitrarily large v , there exists a distribution with optimal function in F_0 such that $\liminf_{n \rightarrow \infty} \Pr(k^* \geq v) > 0$.*

Now, let's consider a distribution that depends on n . This is interesting, as it provides lower bounds on what can be guaranteed, for a given value of n , independently of the distribution. For technical reasons, and without loss of generality within renumbering of the F_k provided that the VC-dimension goes to infinity as $k \rightarrow \infty$, we assume that F_{v+1} has a VC-dimension larger than F_v .

Then:

- Consider n_0 such that for any $n \geq n_0$, there exists v such that $p(v, n) \geq 1/n$. Choose v_n , the maximal value of v such that $p(v, n) \geq 1/n$ and F_v has VC-dimension greater than the VC-dimension of F_{v-1} .
- Define A a set shattered by F_v and not by F_{v-1} . Consider the uniform distribution on $A \times \{0, 1\}$.

- v_n and the distribution are well-defined, if $n \in S$, where $S = [[n_0, \infty[[$;
- $\lim_{n \rightarrow \infty} v_n = \infty$ as for any v , $p(v, n) \rightarrow 1/4^V$ as $n \rightarrow \infty$;
- Consider $N = |\{n \in S; k^* \geq v\}|$;
- N has infinite expectation $E N \geq \sum_{n \in S} 1/n$;
- Therefore, infinitely often (almost surely), $k^* \geq v_n$ and $v_n \rightarrow \infty$, therefore $\limsup k^* = \infty$.

We have therefore shown, *with a distribution dependent on n* , that

$$\limsup_{n \rightarrow \infty} k^* \rightarrow \infty.$$

This leads to this other negative theorem about the control of bloat by hold out.

Proposition 4.3 (Bloat can not be controlled by hold out with pairing, second result). $\limsup_n k^* = \infty$, *where the distribution depends on n but is always such that an optimal function lies in F_0 .*

This result above is in the setting of a distribution which depends on n ; it is of course not interesting for modelizing the evolution of one particular problem as the number of examples increases, but it shows that no bound on $k^*(n)$ for $n \geq n_0$ can be provided, whatever may be n_0 , for hold out with pairing, unless the distribution of problems is taken into account.

4.1.3 Cross-Validation for the Control of Bloat

We now turn our attention to the case of cross-validation. We formalize N-folds-cross-validation as follows:

$$\begin{aligned} f_k^i &= \arg \min_{F_k} L(\cdot, X_k^i) \\ X_k^i &= (X_k^1, X_k^2, \dots, X_k^{i-1}, X_k^{i+1}, X_k^{i+2}, \dots, X_k^N) \text{ for } i \leq N \\ k^* &= \arg \min \frac{1}{N} \sum_{i=1}^N L(f_k^i, X_k^i) \end{aligned}$$

where for any i and k , X_k^i is a sample of n points.

Greedy cross-validation could be considered as in the case of hold out above: all X_k^i could be independent. This leads to the same result (for some distribution, $k^* \rightarrow \infty$) with roughly the same proof. We therefore only consider cross-validation with pairing, i.e. $\forall i, k, k', X_k^i = X_{k'}^i$. For short, we note $X_k^i = X^i$.

Cross-validation is usually applied with finite fixed values of N . People often consider that N should be between 3 and 10, closer to 3 when n increases to infinity. We therefore consider this case; note that the case $N(n) \rightarrow \infty$ is not covered by our results; this could be an interesting further work.

We consider cross-validation for computing k^* . We note \hat{P}_i the empirical law associated to X^i . We consider A a set of points shattered by F_v , $|A| = V$, A not shattered by F_{v-1} . We assume that the distribution of examples is, for x , uniform on A , and for y , independently of x , uniform on $\{0, 1\}$.

We consider $f \in F_v$ realizing a dichotomy of A that is not realized by F_{v-1} . We define E_i the event $\{\forall a \in A; \hat{P}_i(x = a) = 0 \text{ or } \hat{P}_i(y = f(a)|x = a) > \frac{1}{2}\}$.

The probability of E_i goes to $(\frac{1}{2})^{|A|}$ as $n \rightarrow \infty$. The probability of $E = \cap_i E_i$ goes to $(\frac{1}{2})^{N|A|} > 0$ as $n \rightarrow \infty$. In particular, almost surely, infinitely often as $n \rightarrow \infty$, E occurs. When E occurs,

- all the f_k^i are equal to f provided that $k \geq v$;
- for any $g \in F_{v-1}$, $L(g, X_k^i) > L(f, X_k^i)$;
- therefore, $k^* \geq v$.

As E is lower bounded by some positive probability for any value of n , E occurs infinitely often. We therefore have the following result, summarizing this elements.

Theorem 4.4. *Assume that F_k has a VC-dimension going to ∞ as $k \rightarrow \infty$. One can not avoid bloat with only hold-out or cross-validation, in the sense that*

with paired hold out, or greedy hold-out, or cross-validation, for any V , there exists some distribution for which almost surely, $k^ > V$ infinitely often whereas an optimal function lies in F_0 .*

Note that some propositions above show in some cases stronger forms of bloat.

If we consider greedy hold-out, hold out with pairing and cross-validation with pairing, then:

- For some well-chosen distribution of examples, greedy hold-out almost surely leads to (i) $k^* \rightarrow \infty$ if $N \rightarrow \infty$ (ii) k^* asymptotically uniformly distributed in $[[0, N]]$ if N finite, whereas an optimal function lies in F_0 (theorem 4.1).
- Whatever may be $V = VCdim(F_v)$, for some well-chosen distribution, hold-out with pairing almost surely leads to $k^* > V$ infinitely often whereas an optimal function lies in F_0 (proposition 4.2).
- Whatever may be $V = VCdim(F_v)$, for some well-chosen distribution, cross-validation with pairing almost surely leads to $k^* > V$ infinitely often whereas an optimal function lies in F_0 .

4.2 Time-Complexity Requirement in the Fitness

We have shown UC results in previous sections, with convergence rates $1/\sqrt{n}$ provided that both the programs used in the genetic programming tool and the target concept have bounded complexity. This complexity was measured in terms of VC-dimension, through bounds on the number of times each instruction lies in the code, but also through bounds on the number of times these instructions are run. This implies that we deal with programs with bounded time-complexity. We here show that such an hypothesis is necessary, at least with the set of instructions that we have considered.

Consider any learning algorithm working on a sequence of i.i.d. examples $(x_1, y_1), \dots, (x_n, y_n)$ and outputting a program. We formalize as follows the asymptotic convergence rate of an algorithm that does not forbid slow programs.

Definition 4.5 (Convergence rate of a learning program that allows slow programs). *A learning program has convergence rate $(a_n)_{n \in \mathbb{N}}$, if a_n decreases to 0 and if for any distribution \Pr such that there exists some (always halting) function f such that $\Pr(f(x) = y) = 1$, the function f_n provided by the learning program working on n examples has expected error rate $E\Pr(f_n(x) \neq y) = O(a_n)$.*

We show in the sequel that indeed, there exists no convergence rate of learning programs that allow slow programs. In particular, convergence rates as shown in Theorem 3.1, i.e. a guaranteed convergence rate in $O(1/\sqrt{n})$ when an optimal function has bounded complexity (including time!), can not occur. In the sequel, we assume that the reader is familiar with statistical learning theory and shattering properties; the interested reader is referred to [8].

We will need results about learning on families of functions that shatter infinite sets. First, we show that short programs (on real numbers) can generate functions with infinite VC-dimension.

Lemma 4.6 (Parameterizable program that shatters an infinite set). *The following program has bounded length, only one parameter α , but generates as $\alpha \in [0, 1]$ a family of functions which shatters an infinite set:*

- Consider x the entry in $]0, 1[$ and $\alpha \in [0, 1]$ a parameter;
- Label *BEGIN*.
- If $x > \frac{1}{2}$, go to *FINISH*.
- $x \leftarrow 2x$.
- $\alpha \leftarrow 2\alpha$.

- *Goto BEGIN.*
- *Label FINISH.*
- *If $\alpha > 1$ then $\alpha \leftarrow \alpha - 1$; goto FINISH.*
- *If $\alpha \geq 0.5$, output 1 and stop.*
- *Output 0 and stop.*

Proof. Consider $(\alpha, x) \in [0, 1] \times [0, 1]$ and their binary representations $x = \sum_{i \geq 0} x_i 2^{-i}$, $\alpha = \sum_{i \geq 0} \alpha_i 2^{-i}$. This program shifts α and x to the left until $x > \frac{1}{2}$, i.e. the program applies $j - 1$ shifts with j minimal such that $x_j = 1$. It then replies 1 if and only if α , after shift, has its digit α_j equal to 1. Therefore, this program can realize any dichotomy of $\{\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots\}$. This is exactly the definition of the fact that this set is shattered. So, we have shown that a family of functions shattering an infinite set was included in the set of programs with bounded length. ■

We now have all required tools for proving the following.

Theorem 4.7 (No convergence rate when working on slow programs). *What ever may be the sequence a_1, \dots, a_n, \dots decreasing to 0, for any learning program, for some f which always halts and with bounded length, the expectation of $\Pr(P_n(x) \neq y)$ is not $O(a_n)$.*

Proof. This is an immediate consequence of the lemma above, and of the classical lower bound on learning capabilities of algorithms working on spaces of functions shattering an infinite set. See e.g. [8, chap. 14]. ■

Some complements (non-continuous domains, time-complexity of the optimization with some penalization terms) can be found in [26].

5 Experimental Results

Some theoretical elements presented in Sections 2 and 3 are verified experimentally in this section. The experimentation are conducted using Koza-style GP [13], with a problem setup similar to the classical symbolic regression example, modified for binary classification. This is a simpler case than the computing machine presented in Lemma 2.2, although the conclusions previously made still apply.

The GP branches used are the addition, subtraction, multiplication, protected division, and if-less-than. This last branch takes four arguments, returning the third argument if the first argument is less than the second one, otherwise returning the fourth argument. The GP terminals are the x variable, and the 0 and 1 constants. The learning task consists in minimizing the error $e(i)$ between the desired output $y_i = \{-1, 1\}$ and the obtained output \hat{y}_i of the tested GP tree for the x_i input, as in the following:

$$e(i) = \max(1 - y_i \hat{y}_i, 0)$$

The fitness measure used in the experiments consists in minimizing the sum of the errors to which is added a complexity factor that approximate the VC-dimension of the GP program:

$$f = \frac{1}{s} \sum_{i=1}^s e(i) + k \sqrt{\frac{t^2 \log_2(t)}{s}},$$

where t is the number of nodes of the GP program tested, s is the number of test cases used for fitness evaluation, and k is a trade-off weight in the composition of the complexity penalization relatively to the accuracy term.

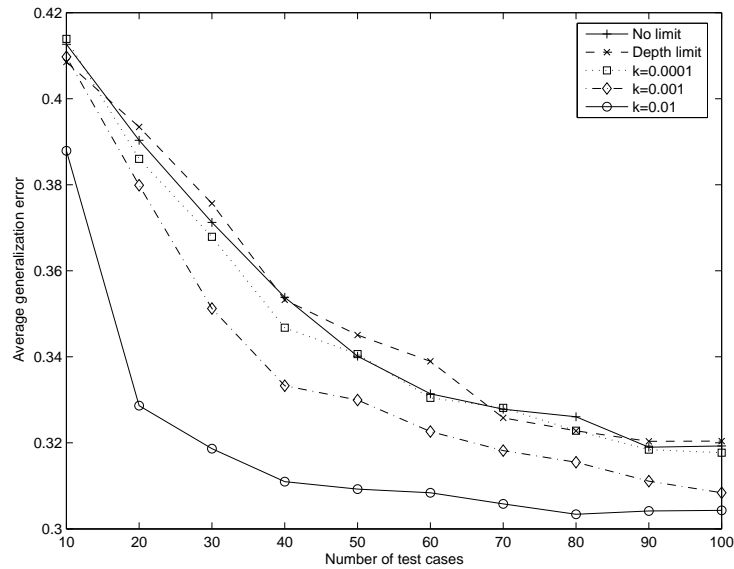
The s test cases are distributed uniformly in $x_i \in [0, 1]$, with associated $y_i = \{-1, 1\}$. For $x_i < 0.4$, each y_i are equal to 1 with probability 0.25 (so $y_i = -1$ with probability 0.75), for $x_i \in [0.4, 0.6[$, $y_i = 1$ with probability 0.5, and for $x_i \geq 0.6$, $y_i = 1$ with probability 0.75. Thus, the associated

classifier with best generalization capabilities would return $y_i^* = -1$ for $x_i < 0.4$, $y_i^* = 1$ for $x_i \geq 0.6$ and a random output for $x_i \in [0.4, 0.6[$, with a minimal generalization error of 0.3. After the evolutions, each best-of-run classifier is thus evaluated by a fine sampling of the input space, with the generalization error evaluated as the difference between the output given by the tested best-of-run classifier and the output obtained by a classifier with best generalization capabilities.

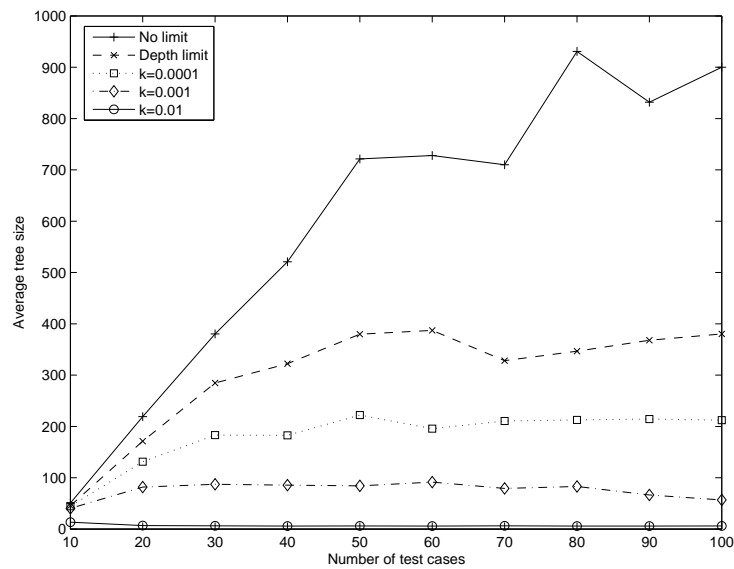
Five types of GP evolutions have been tested: i) no limitation on the tree size (no depth limit and complexity trade-off $k = 0$), ii) depth limitation on the tree size of 17 levels (complexity trade-off $k = 0$), iii) soft complexity penalty in the fitness (complexity trade-off $k = 0.0001$), iv) medium complexity penalty in the fitness (complexity trade-off $k = 0.001$), and v) important complexity penalty in the fitness (complexity trade-off $k = 0.01$). For the three last approaches, the depth limitation of 17 levels is still maintained. The selection method used is lexicographic parsimony pressure [18], that is regular tournament selection 4 participants, with the smallest participant taken in case of ties. Other GP parameters are: population of 1000 individuals; evolutions on 200 generations; crossover probability of 0.8; subtree, swap and shrink mutation of probability 0.05 each; and finally half-and-half initialization with maximal depth of 5. All the experiments have been implemented using the GP facilities of the Open BEAGLE² [10] C++ framework for evolutionary computations.

The experiments have been conducted different number of test cases varying from $s = 10$ to $s = 100$ by steps of 10. One hundred evolutions is done for each combinations of approaches tested and number of test cases, for a total of 50000 evolutions. Figure 4 shows the average generalization errors and tree size obtained for the different approaches in function of the number of test cases used for fitness evaluation.

²<http://beagle.gel.ulaval.ca>



(a)



(b)

Figure 4: Generalization errors and tree sizes observed for different size limitations. Figure (a) shows the average generalization errors observed, with apparently better results for the approaches where the fitness includes some parsimony pressure. Figure (b) shows the average tree sizes obtained, where important bloat is observed for the no limitation and maximum depth limitations

These results show that bloat occurs when no limitation of size occurs, even when lexicographic parsimony pressure is used (see curve *No limit* of Figure 4b), which validates Theorem 2.6. Then, as stated by Theorem 2.4, UC is achieved using moderate complexity penalization in the fitness measure, with a convergence toward optimal generalization error of 0.3 (see curve $k=0.001$ of Figure 4a). Third, as predicted by Theorem 3.1, increasing the penalization leads to both UC and no bloat (see curve $k=0.01$ of both Figures 4a and 4b). Note that Theorem 2.6 asserts that this result cannot be achieved by *a priori* scaling of the complexity, and that Section 4 shows that this can not be achieved by cross-validation.

6 Conclusion

In this paper, we have proposed a theoretical study of two important issues in Genetic Programming (GP) known as Universal Consistency (UC) and code bloat. We have shown that the understanding of the bloat phenomenon in GP could benefit from classical results from statistical learning theory.

The first limit of our work is the fact that all these results consider that GP finds a program which is empirically the best, in the sense that given a set of test cases and a fitness function based on the empirical error (and possibly including some parsimony penalization), it will be assumed that GP does find one program in that search space that minimizes this fitness – and it is the behavior of this ideal solution, which is a random function of the number of test cases, that is theoretically studied.

Of course, we all know that GP is not such an ideal search procedure, and hence such results might look rather far away from GP practice, where the user desperately tries to find a program that gives a reasonably low empirical approximation error. Nevertheless, UC is vital for the practitioner too: indeed, it

would be totally pointless to fight to approximate an empirically optimal function without any guarantee that this empirical optimum is anywhere close to the ideal optimal solution we are in fact looking for. Furthermore, the bloat-related results give some useful hints about the type of parsimony that has a chance to efficiently fight the unwanted bloat, while maintaining the UC property. Experiments confirm the results.

Application of theorems from learning theory has led to two original outcomes with both positive and negative results. Firstly, results on UC of GP: there is almost sure asymptotic convergence to the optimal error rate in the context of binary classification with GP. Secondly, results on code bloat: i) if the ideal target function does not have a finite description then code bloat is unavoidable (structural bloat), and ii) code bloat can be avoided by simultaneously bounding the length of the programs with some *ad hoc* limit and using some parsimony pressure in the fitness function (functional bloat). An important point is that all methods leading to no-bloat use a regularization term; in particular, cross-validation or hold-out methods do not reach no-bloat. Also, all methods ensuring convergence rates (with the set of instructions we have chosen) have some form of regularization using time-complexity

Interestingly, all those results (both positive and negative) about code bloat are also valid in different contexts, such as for instance that of Neural Networks (the number of neurons replaces the complexity of GP programs). Moreover, results presented here are not limited to the scope of binary classification problems, and may be applied to variable length representation algorithms in different contexts such as control or identification tasks.

Finally, going back to the debate about the causes of bloat in practice, it is clear that our results can only partly explain the actual cause of bloat in a real GP run – and tend to give arguments to the “fitness causes bloat” explanation

[16]. It might be possible to study the impact of size-preserving mechanisms (e.g. specific variation operators, like size-fair crossover [15] or fair mutations [17]) as somehow contributing to the regularization term in our final result ensuring both UC and no bloat.

Acknowledgements

This work was supported in part by the PASCAL Network of Excellence, and by postdoctoral fellowships from the ERCIM (Europe) and the FQRNT (Québec) to C. Gagné. We thank Bill Langdon for very helpful comments.

A Elements of Statistical Learning Theory

In the frameworks of regression and classification, statistical learning theory [27] is concerned with giving some bounds on the generalization error (i.e. the error on yet unseen data points) in terms of the actual empirical error and some fixed quantity depending only on the search space. More precisely, we will use here the notion of *Vapnik-Chervonenkis dimension* (in short, VC-dimension or *VCdim*) of a space of functions. Roughly, the VC-dimension is an indicator of the theoretical capability of a learning machine to discriminate data. It is often used to in the computation of bounds on the difference between the empirical error and the generalization error.

Consider a set of s examples $(x_i, y_i)_{i \in \{1, \dots, s\}}$. These examples are drawn from a distribution P on the couple (X, Y) . They are independent identically distributed, $Y = \{0, 1\}$ (classification problem), and typically $X = \mathbb{R}^d$ for some dimension d . For any function f , define the *loss* $L(f)$ to be the *expectation* of $|f(X) - Y|$. *Similarly, define the empirical loss* $\hat{L}(f)$ as the loss observed on the examples: $\hat{L}(f) = \frac{1}{s} \sum_i |f(x_i) - y_i|$. Finally, define L^* , the *Bayes error*, as the smallest possible generalization error for any mapping from X to $\{0, 1\}$.

The following four theorems are well-known in the statistical learning community:

Theorem A.1 (Bound on the empirical risk with finite VC-dimension). *Consider \mathcal{F} a family of functions from a domain X to $\{0, 1\}$ and V its VC-dimension. Then, for any $\epsilon > 0$,*

$$\Pr \left(\sup_{P \in \mathcal{F}} |L(P) - \hat{L}(P)| \geq \epsilon \right) \leq 4 \exp(4\epsilon + 4\epsilon^2) s^{2V} \exp(-2s\epsilon^2),$$

and for any $\delta \in]0, 1]$,

$$\Pr \left(\sup_{P \in \mathcal{F}} |L(P) - \hat{L}(P)| \geq \epsilon(s, V, \delta) \right) \leq \delta,$$

where

$$\epsilon(s, V, \delta) = \sqrt{\frac{4 - \log(\delta/(4s^{2V}))}{2s - 4}}.$$

Interpretation. This theorem states that in a family of finite VC-dimension, the empirical errors and the generalization errors are probably closely related.

Proof. See [8, Th. 12.8, p. 206]. ■

Other forms of this theorem have no $\log(n)$ factor; they are known as Alexander's bound, but the constant is so large that this result is not better than the result above unless s is huge (see [8, p. 207]). If $s \geq 64/\epsilon^2$,

$$\Pr \left(\sup_{P \in \mathcal{F}} |L(P) - \hat{L}(P)| \geq \epsilon \right) \leq 16 (\sqrt{s}\epsilon)^{4096V} \exp(-2s\epsilon^2)$$

We classically derive the following result from Theorem A.1:

Theorem A.2 (Convergence of the empirical error to the generalization error with infinite number of examples). *For $s \geq 0$, consider \mathcal{F}_s a family of functions from a domain X to $\{0, 1\}$ and V_s its VC-dimension. Then,*

$$\sup_{P \in \mathcal{F}_s} |L(P) - \hat{L}(P)| \rightarrow 0 \text{ as } s \rightarrow \infty$$

almost surely provided that $V_s = o(s/\log(s))$ (i.e. $V_s \log(s)/s \rightarrow 0$).

Interpretation. The maximal difference between the empirical error and the generalization error goes almost surely to 0 if the VC-dimension is finite.

Proof. We use the classical Borel-Cantelli lemma, for any $\epsilon \in [0, 1]$:

$$\begin{aligned} & \sum_{s \geq 64/\epsilon^2} \Pr(|L(P) - \hat{L}(P)| > \epsilon) \\ & \leq 16 \sum_{s \geq 64/\epsilon^2} (\sqrt{s}\epsilon)^{4096V_s} \exp(-2s\epsilon^2), \\ & = 16 \sum_{s \geq 64/\epsilon^2} \exp(4096V_s(\log(\sqrt{s}) + \log(\epsilon)) - 2s\epsilon^2), \end{aligned}$$

which is finite as soon as $V_s = o(s/\log(s))$. ■

Theorem A.3 (Universal consistency with finite VC-dimension). *Let $\mathcal{F}_1, \dots, \mathcal{F}_k, \dots$ families of functions with finite VC-dimensions, $VCdim(\mathcal{F}_i) = V_i$, and let $\mathcal{F} = \cup_n \mathcal{F}_n$. Then, given s examples, consider $\hat{P}_s \in \mathcal{F}_s$ minimizing the empirical risk \hat{L} among F_s . Then, if $V_s = o(s/\log(s))$ and $V_s \rightarrow \infty$, for any $\delta \in]0, 1]$,*

$$\begin{aligned} \Pr\left(L(\hat{P}_s) \leq \hat{L}(\hat{P}_s) + \epsilon(s, V_s, \delta)\right) & \geq 1 - \delta, \\ \Pr\left(L(\hat{P}_s) \leq \inf_{P \in \mathcal{F}_s} L(P) + 2\epsilon(s, V_s, \delta)\right) & \geq 1 - \delta. \end{aligned}$$

Also, $L(\hat{P}_s) \rightarrow \inf_{P \in \mathcal{F}} L(P)$ almost surely. Note that for a well chosen family of functions (typically programs), $\inf_{P \in \mathcal{F}} L(P) = L^*$ for any distribution. Thus, this theorem leads to universal consistency (i.e. $\forall P \ L(\hat{P}_s) \rightarrow L^*$), for a well-chosen family of functions.

Interpretation. If the VC-dimension increases slowly enough as a function of the number of examples, then the generalization error goes to the optimal one. If the family of functions is well-chosen, this slow increase of VC-dimension leads to universal consistency.

Proof. See [8, Th. 18.2, p. 290] and [11]. ■

In the following theorem, we use d', t', q' instead of d, t, q for the sake of consistency of notations as in other results we need d, t and q but we will apply Theorem A.4 with some *ad hoc* d', t' and q' .

Theorem A.4 (Bound on VC-dimension for a computing machine). *Let $H_{h,d'} = \{x \mapsto h(a, x); a \in R^{d'}\}$ where h can be computed with at most t' operations among $\alpha \mapsto \exp(\alpha); +, -, \times, /$; jumps conditioned on $>, \geq, <, \leq, =$; output 0; output 1 (same set of instructions as in other parts of this paper). We note $H_{h,d'}$ as H when there is no ambiguity.*

Then: $VCdim(H) \leq t'^2 d' (d' + 19 \log_2(9d'))$.

Furthermore, if $\exp(\cdot)$ is used at most q' times, and if there are at most t' operations executed among arithmetic operators, conditional jumps, exponentials, then:

$$\pi(H, m) \leq 2^{(d'(q'+1))^{2/2}} (9d'(q'+1)2^t)^{5d'(q'+1)} \left(em(2^{t'} - 2) / d' \right)^{d'}$$

where $\pi(H, m)$ is the m^{th} shattering coefficient of H , and hence

$$VCdim(H) \leq (d'(q'+1))^2 + 11d'(q'+1)(t' + \log_2(9d'(q'+1)))$$

Finally, if $q' = 0$ then $VCdim(H) \leq 4d'(t' + 2)$.

Interpretation. The VC-dimension of the set of the possible parameterizations of a program as defined above is bounded.

Proof. See [1, 8.14 and 8.4]. ■

Theorem A.5 (Structural risk minimization). *Let $\mathcal{F}_1, \dots, \mathcal{F}_k \dots$ with finite VC-dimensions $VCdim(\mathcal{F}_i) = V_i$. Let $\mathcal{F} = \cup_n \mathcal{F}_n$. Assume that all distributions lead to $L_{\mathcal{F}} = L^*$ where L^* is the optimal possible error (spaces of functions ensuring this exist). Given s examples, consider $f \in \mathcal{F}$ minimizing $\hat{L}(f) + \sqrt{\frac{32}{s} V(f) \log(es)}$, where $V(f)$ is V_k with k minimal such that $f \in \mathcal{F}_k$. Then:*

- If additionally one optimal function belongs to \mathcal{F}_k , then for any s and ϵ such that $V_k \log(es) \leq s\epsilon^2/512$, the generalization error is larger than ϵ with probability at most $\Delta \exp(-s\epsilon^2/128) + 8s^{V_k} \exp(-s\epsilon^2/512)$ where $\Delta = \sum_{j=1}^{\infty} \exp(-V_j)$ is assumed finite;
- The generalization error, with probability 1, converges to L^* .

Interpretation. The optimization of a compromise between empirical accuracy and regularization leads to the same properties as in Theorem A.3, plus a stronger convergence rate property.

Proof. See [27] and [8, p. 294]. ■

References

- [1] M. Antony and P.L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999.
- [2] W. Banzhaf and W. B. Langdon. Some considerations on the reason for bloat. *Genetic Programming and Evolvable Machines*, 3(1):81–91, 2002.
- [3] W. Banzhaf, P. Nordin, R.E. Keller, and F.D. Francone. *Genetic Programming : an introduction*. Morgan Kaufmann Publisher Inc., San Francisco, CA, USA, 1998.
- [4] Stefan Bleuler, Martin Brack, Lothar Thiele, and Eckart Zitzler. Multiobjective genetic programming: Reducing bloat using SPEA2. In *Proceedings of the 2001 Congress on Evolutionary Computation CEC2001*, pages 536–543, COEX, World Trade Center, 159 Samseong-dong, Gangnam-gu, Seoul, Korea, 27-30 2001. IEEE Press.
- [5] T. Blickle and L. Thiele. Genetic programming and redundancy. In J. Hopf, editor, *Genetic Algorithms Workshop at KI-94*, pages 33–38. Max-Planck-Institut für Informatik, 1994.
- [6] J. M. Daida, R. R. Bertram, S. A. Stanhope, J. C. Khoo, S. A. Chaudhary, O. A. Chaudhri, and J. A. Ii Polito. What makes a problem GP-Hard? Analysis of a tunably difficult problem in genetic programming. *Genetic Programming and Evolvable Machines*, 2(2):165 – 191, 2001.
- [7] Edwin D. De Jong, Richard A. Watson, and Jordan B. Pollack. Reducing bloat and promoting diversity using multi-objective methods. In *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO-2001*, pages 11–18, San Francisco, CA, 2001. Morgan Kaufmann Publishers.

- [8] L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, 1997.
- [9] A. Ekart and S. Nemeth. Maintaining the diversity of genetic programs. In *EuroGP '02: Proceedings of the 5th European Conference on Genetic Programming*, pages 162–171, London, UK, 2002. Springer-Verlag.
- [10] C. Gagné and M. Parizeau. Genericity in evolutionary computation software tools: Principles and case study. *International Journal on Artificial Intelligence Tools*, 15(2):173–194, April 2006.
- [11] U. Grenander. *Abstract Inference*. Wiley, New York, 1981.
- [12] S. Gustafson, A. Ekart, E. Burke, and G. Kendall. Problem difficulty and code growth in genetic programming. *Genetic Programming and Evolvable Machines*, 4(3):271–290, 2004.
- [13] J. R. Koza. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA, USA, 1992.
- [14] W. B. Langdon. The evolution of size in variable length representations. In *IEEE International Congress on Evolutionary Computations (ICEC 1998)*, pages 633–638. IEEE Press, 1998.
- [15] W. B. Langdon. Size fair and homologous tree genetic programming crossovers. *Genetic Programming And Evolvable Machines*, 1(1/2):95–119, 2000.
- [16] W. B. Langdon and R. Poli. Fitness causes bloat: Mutation. In *Late Breaking Papers at GP'97*, pages 132–140. Stanford Bookstore, 1997.
- [17] W. B. Langdon, T. Soule, R. Poli, and J. A. Foster. The evolution of size and shape. In *Advances in Genetic Programming III*, pages 163–190. MIT Press, 1999.

- [18] S. Luke and L. Panait. Lexicographic parsimony pressure. In *GECCO 2002: Proceedings of the Genetic and Evolutionary Computation Conference*, pages 829–836. Morgan Kaufmann Publishers, 2002.
- [19] N. F. McPhee and J. D. Miller. Accurate replication in genetic programming. In *Genetic Algorithms: Proceedings of the Sixth International Conference (ICGA95)*, pages 303–309, Pittsburgh, PA, USA, 1995. Morgan Kaufmann.
- [20] P. Nordin and W. Banzhaf. Complexity compression and evolution. In *Genetic Algorithms: Proceedings of the Sixth International Conference (ICGA95)*, pages 310–317, Pittsburgh, PA, USA, 15-19 1995. Morgan Kaufmann.
- [21] A. Ratle and M. Sebag. Avoiding the bloat with probabilistic grammar-guided genetic programming. In *Artificial Evolution VI*. Springer Verlag, 2001.
- [22] S. Silva and J. Almeida. Dynamic maximum tree depth : A simple technique for avoiding bloat in tree-based GP. In *Genetic and Evolutionary Computation – GECCO-2003*, LNCS, pages 1776–1787. Springer-Verlag, 2003.
- [23] Sara Silva and Ernesto Costa. Dynamic limits for bloat control: Variations on size and depth. In *GECCO (2)*, pages 666–677, 2004.
- [24] T. Soule. Exons and code growth in genetic programming. In *European Conference on Genetic Programming (EuroGP 2002)*, volume 2278 of *LNCS*, pages 142–151. Springer-Verlag, 2002.

- [25] T. Soule and J. A. Foster. Effects of code growth and parsimony pressure on populations in genetic programming. *Evolutionary Computation*, 6(4):293–309, 1998.
- [26] O. Teytaud. Why simulation-based approaches with combined fitness are a good approach for mining spaces of turing-equivalent functions. In *Proc. of the IEEE Congress on Evolutionary Computation (CEC 2006)*, 2006.
- [27] V. Vapnik. *The Nature of Statistical Learning Theory*. Springer, 1995.
- [28] B.-T. Zhang and H. Mühlenbein. Balancing accuracy and parsimony in genetic programming. *Evolutionary Computation*, 3(1), 1995.
- [29] B.-T. Zhang, P. Ohm, and H. Mühlenbein. Evolutionary induction of sparse neural trees. *Evolutionary Computation*, 5(2):213–236, 1997.