



General-purpose Hierarchical Optimisation of Machine Learning Pipelines with Grammatical Evolution

Suilan Estevez-Velarde, Yoan Gutiérrez, Yudiivián Almeida-Cruz, Andrés Montoyo

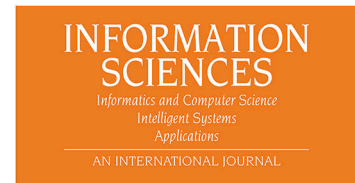
PII: S0020-0255(20)30698-8
DOI: <https://doi.org/10.1016/j.ins.2020.07.035>
Reference: INS 15673

To appear in: *Information Sciences*

Received Date: 10 January 2020
Accepted Date: 10 July 2020

Please cite this article as: S. Estevez-Velarde, Y. Gutiérrez, Y. Almeida-Cruz, A. Montoyo, General-purpose Hierarchical Optimisation of Machine Learning Pipelines with Grammatical Evolution, *Information Sciences* (2020), doi: <https://doi.org/10.1016/j.ins.2020.07.035>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



General-purpose Hierarchical Optimisation of Machine Learning Pipelines with Grammatical Evolution

Suilan Estevez-Velarde^{a,*}, Yoan Gutiérrez^{b,c}, Yudiivián Almeida-Cruz^a,
Andrés Montoyo^{b,c}

^a*School of Math and Computer Science, University of Havana, Cuba*

^b*University Institute for Computing Research (IUII), University of Alicante, Spain*

^c*Department of Software and Computing Systems, University of Alicante, Spain*

Abstract

This paper introduces Hierarchical Machine Learning Optimisation (HML-Opt), an AutoML framework that is based on probabilistic grammatical evolution. HML-Opt has been designed to provide a flexible framework where a researcher can define the space of possible pipelines to solve a specific machine learning problem, which can range from high-level decisions about representation and features to low-level hyper-parameter values. The evaluation of HML-Opt is presented via two different case studies, both of which demonstrate that it is competitive with existing AutoML tools on a variety of benchmarks. Furthermore, HML-Opt can be applied to novel problems, such as knowledge extraction from natural language text, whereas other techniques are insufficiently flexible to capture the complexity of these scenarios. The source code for HML-Opt is available online for the research community.

Keywords: AutoML, Grammatical evolution, Evolutionary computation, Supervised learning, Natural language processing

1. Introduction

The research process for solving a machine learning problem often involves experimenting with several different approaches (neural networks, supervised

*Corresponding author.

Email addresses: `sestevez@matcom.uh.cu` (Suilan Estevez-Velarde),
`ygutierrez@dlsi.ua.es` (Yoan Gutiérrez), `yudy@matcom.uh.cu` (Yudiivián
Almeida-Cruz), `montoyo@dlsi.ua.es` (Andrés Montoyo)

4 classifiers, clustering algorithms, etc.). Designing an effective solution to any
 5 one of these problems often involves a large experimentation phase where re-
 6 searchers use different datasets, algorithms, and specific parameters. As an
 7 example, Figure 1 shows a hypothetical pipeline composed of several steps.
 8 In each step, different options are available. Suitable combinations of these
 9 options yield different values for the performance metric that is being eval-
 10 uated. These hypothetical steps can range from applying some data prepro-
 11 cessing techniques to selecting specific algorithms and further determining
 12 the values of their hyperparameters.

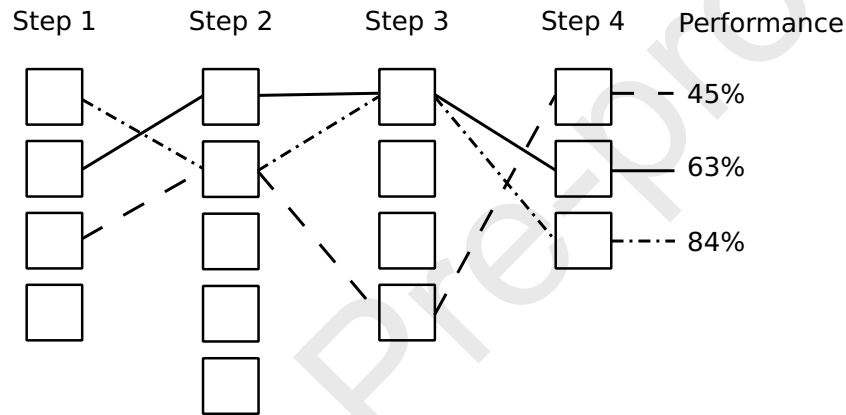


Figure 1: Illustrative example of an abstract pipeline, in which several steps can be performed in different ways, each combination obtaining a different performance (e.g, classification accuracy, recall, etc.)

13 Exploring all the possible combinations of algorithms and parameters for
 14 a given problem can be unfeasible, since the number of possibilities is often
 15 exponential with respect to the number of steps in a pipeline. Furthermore,
 16 when some algorithms have numerical (discrete or continuous) parameters
 17 (e.g., regularisation rate, number of neurons in a neural network layer), it is
 18 impossible to evaluate all the combinations. This problem is further com-
 19 plicated by the fact that each experiment may have a high computational
 20 complexity (e.g., training from scratch a neural network on a large dataset).
 21 Researchers are often compelled to select a small number of possibilities based
 22 on previous experience and domain knowledge about the problem at hand.

23 The automation of this lengthy experimentation process is denominated
 24 Automatic Machine Learning (AutoML). AutoML is an increasingly growing
 25 field that has been applied for finding optimal machine learning pipelines in a

26 variety of scenarios. As an example, in computer vision, where several neural
27 network architectures have been extensively explored, Zoph et al. [44] applies
28 reinforced learning to learn to build the optimal neural network architecture
29 for any given image dataset. Likewise, general frameworks such as Auto-
30 Sklearn [12], Auto-Weka [39], or Auto-Keras [15] have appeared, based upon
31 existing machine learning frameworks, that automatically explore different
32 combinations of algorithms available in such frameworks.

33 Current AutoML techniques focus mostly on a specific subset of algo-
34 rithms, often tailored to a specific framework or tool-set. Solving complex
35 problems, on the other hand, requires the combination of different tools
36 that might not be available in a single framework. Besides, challenging ma-
37 chine learning problems are not restricted to finding the best architecture
38 or hyper-parameter set for a given algorithm, but often involve higher-level
39 decisions. For example, in natural language processing, prior to deciding
40 which machine learning algorithm to use, the researcher must decide on rep-
41 resentation (e.g., whether to use embeddings, and which ones specifically),
42 features (e.g., whether to incorporate knowledge-based features), preprocess-
43 ing (e.g., whether to remove stop-words or to apply stemming), etc. Advanc-
44 ing towards the automation of this process requires more expressive tools
45 that allow a researcher to define a complex and problem-specific space of
46 possible pipelines including everything from high-level decisions to low-level
47 hyper-parameter ranges.

48 The main objective of this research is to present Hierarchical Machine
49 Learning Optimisation (HML-Opt), an AutoML technique based on proba-
50 bilistic grammatical evolution. This work extends and generalises previous
51 research by Estevez-Velarde et al. [11]. The most important contributions of
52 this research are:

- 53 • HML-Opt allows a researcher to define a large and complex space of
54 possible pipelines to solve a specific problem, potentially combining
55 different frameworks and technologies.
- 56 • HML-Opt automatically searches this space of pipelines, with given
57 time frames and computational resources, to find an optimal or close
58 to optimal architecture.
- 59 • Meaningful statistics and insights about the problem are extracted from
60 the experimentation process.

- 61 • Freely available source code is provided for the research community¹.

62 This rest of the paper is organised as follows. Section 2 presents a brief
 63 review of the relevant concepts and related work. Section 3 introduces HML-
 64 Opt, describing its overall design and relevant implementation details. Sec-
 65 tion 4 presents two case studies involving different datasets, ranging from
 66 simple classification problems to state-of-the-art knowledge discovery prob-
 67 lems. Finally, Section 5 discusses the main contributions and highlights of
 68 the research and the case study results, and Section 6 presents the main
 69 conclusions and future lines of research.

70 2. Related work

71 The idea of designing meta-algorithms to select the best algorithms for
 72 specific problem domains is a recurrent trend in artificial intelligence research,
 73 which has been motivated by several factors. These include: the complex-
 74 ity of parameter tuning in practical problems; the wide variety of existing
 75 algorithms with similar performance; and, the existence of no-free-lunch re-
 76 sults. In the domains of continuous and combinatorial optimisation, hybrid
 77 strategies have been developed to select from different search techniques (e.g.,
 78 different metaheuristics) for a specific problem [4]. Hyperheuristics [3] are
 79 high-level methodologies to select or generate strategies from an underlying
 80 space of heuristic components to solve a specific optimisation problem or class
 81 of problems. A similar formulation is developed in the field of automatic ma-
 82 chine learning (AutoML), in which the underlying components are learning
 83 algorithms and the problem to solve is a supervised learning problem [14].

84 In the context of AutoML, several approaches attempt to find the best
 85 architecture for a given problem either through searching [30] or through
 86 meta-learning [23]. Different approaches exist for AutoML. Among the most
 87 relevant, we can cite: Bayesian optimisation, e.g., Auto-Sklearn [12], Auto-
 88 Weka [39], and Hyperot [22]; Neural Architecture Search (NAS), e.g., Auto-
 89 Keras [15], Zoph and Le [43], Zoph et al. [45]; Genetic Programming, e.g.,
 90 Recipe [9], TPOT [31]; Hierarchical Task Networks, e.g., ML-Plan [29]; and,
 91 other metaheuristics applied to specific problem domains [35].

92 Several metaheuristics have been used as optimisation strategies for in-
 93 telligently searching the best combination of algorithms for a given prob-

¹<https://github.com/knowledge-learning/hp-optimization>

94 lem [19, 18]. Among the most used strategies are genetic algorithms [24] and
 95 evolutionary algorithms [35]. In this last group, grammatical evolution [10]
 96 has been employed for optimising the architecture of neural networks [5] as
 97 well as for directly finding the optimal weights [41, 37, 1]. This increases
 98 the complexity of the optimisation, because the search space is much larger.
 99 Furthermore, there are many efficient optimisation strategies specifically for
 100 tuning neural network weights using variants of gradient descent, such as
 101 SGD [2], RMSProp [40], and Adam [21]. For these reasons, our proposal
 102 attempts to optimise the architecture and high-level hyperparameters while
 103 the actual model parameters are tuned via specialised learning algorithms.
 104 Moreover, our proposal not only applies to neural networks, but in general to
 105 any machine learning pipeline that can be designed as a hierarchical process.

106 Hierarchical representations have shown to be effective for efficiently ex-
 107 ploring this large space [26]. These representations impose a structure in
 108 the search space that helps to guide the optimisation by ensuring that se-
 109 mantically similar elements are closer in the search space. Since grammars
 110 provide an effective and simple method for researchers to define the hierar-
 111 chical structure of a given space [9], our proposal uses grammatical evolution
 112 to impose a semantic structure. The majority of metaheuristic algorithms
 113 deal either with continuous or with discrete values independently because
 114 comparing solutions with a mixture of continuous and discrete values is
 115 harder [17, 42, 16]. HML-Opt has been designed to handle both continu-
 116 ous and discrete values, and factors in that solutions represent hierarchical
 117 pipelines. Therefore, differences between solutions at higher levels of the
 118 hierarchy have greater significance.

119 Table 1 shows the characteristics that define our proposal compared to
 120 some relevant alternatives in the literature. We focus fundamentally on Au-
 121 toML approaches —row 4—, although we include a reference to Grammatical
 122 Evolution applied to probabilistic context-free grammars [20] because it is
 123 similar to our proposal in the theoretical formulation —first column—. We
 124 make a distinction of which techniques can deal with discrete and/or continu-
 125 ous hyperparameters —rows 5 and 6—. Also, we are interested in techniques
 126 that explicitly build a probabilistic description of the search space —row
 127 3—. Finally, we highlight techniques that can be generalised across many
 128 machine learning problems over those that are specifically designed for one
 129 problem, such as supervised classification, or one learning paradigm, such as
 130 neural networks —row 7—.

131 Most of the AutoML approaches presented in Table 2 are black-box, in

	Features	Prob. GE	Auto-Sklearn	Auto-Weka	Real et al.	Recipe	TPOT	Auto-Keras	ML-Plan	Hyperopt	HML-Opt
1	Hierarchical	✓				✓			✓		✓
2	Metaheuristics	✓			✓	✓	✓				✓
3	Space description	✓	✓	✓						✓	✓
4	AutoML		✓	✓	✓	✓	✓	✓	✓	✓	✓
5	Continuous		✓	✓		?		✓		✓	✓
6	Discrete	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
7	General-purpose	✓				✓		✓	✓	✓	✓
8	Year	2015	2015	2013	2017	2017	2018	2016	2018	2014	2019

Table 1: The characteristics of our proposal compared to related work. A question mark (“?”) indicates that we could not find evidence for that characteristic in the corresponding paper.

132 the sense that they define a single generic machine learning pipeline which
 133 is parameterized by a fixed subset of algorithms and their hyperparameters.
 134 These approaches are difficult to apply to problems which are not directly
 135 framed as classic supervised learning from an input feature matrix (eg., NLP
 136 problems where text preprocessing must also be taken into consideration).
 137 A notable exception is the case of Hyperopt [22], which provides a general-
 138 purpose framework that can describe many different pipelines, and the au-
 139 thors provide one pipeline, specifically tailored for `scikit-learn` algorithms.
 140 Likewise, in the case of Recipe [9], the authors predefine different grammars
 141 with fixed sets of algorithms, making the proposal easy to adapt to other
 142 scenarios. A recent technique is ML-Plan [29], which proposes the use of Hi-
 143 erarchical Task Networks as a means for representing the search space, and
 144 two specific implementations, based either on `weka` or `scikit-learn`. HML-
 145 Opt is also designed as an extensible framework that can be tailored to any
 146 specific subset of algorithms, such as `sklearn` or `keras`. However, in con-
 147 trast with other approaches, HML-Opt takes this extensibility a step further,
 148 by giving the user the possibility to completely define the space of possible
 149 pipelines in an intuitive, declarative language, i.e., with a context-free gram-
 150 mar. We illustrate this extensibility in Section 4.2 with an application to
 151 a domain-specific machine learning problem with key restrictions about the
 152 possible pipelines.

153 In terms of representation, the most relevant related works are Recipe [9]

154 and ML-Plan [29] which also define hierarchical search spaces. More specifi-
 155 cally, Recipe also employs a grammar to define the set of possible pipelines.
 156 One key difference in our work is the use of probabilistic grammatical evolu-
 157 tion, as opposed to genetic programming in Recipe, which provides a prob-
 158 abilistic model of the space of pipelines. In ML-Plan the search algorithm
 159 is a modified variant of Monte Carlo tree search designed to deal with the
 160 large computational cost of evaluating machine learning pipelines. Other
 161 AutoML proposals based on Bayesian optimisation (e.g., Auto-Sklearn and
 162 Auto-Weka) also define a probabilistic model of the pipelines, even though it
 163 is not explicitly used to gather insights of the optimisation process. HML-Opt
 164 leverages the use of grammars for defining a flexible and adaptable structure
 165 of pipelines, and the use of a probabilistic model which enables analysis and
 166 understanding of the optimisation process for a specific machine learning
 167 problem.

168 3. Hierarchical Machine Learning Optimisation

169 This section presents Hierarchical Machine Learning Optimisation (HML-
 170 Opt), a technique that searches for the optimal pipeline for a specific machine
 171 learning problem. The core of our proposal is a process that involves three
 172 different parts. First, the researcher designs a grammar that defines the space
 173 of all possible pipelines that are interesting to analyze for a given machine
 174 learning problem. (see Section 3.1). This grammar allows the generation
 175 of different pipelines which can be evaluated for the problem at hand (see
 176 Section 3.2). Finally, an iterative optimisation algorithm continuously gen-
 177 erates and evaluates pipelines searching for the optimal combination (see
 178 Section 3.3). Figure 2 presents a high-level representation of this process.

179 In HML-Opt, the grammar that defines the search space is a probabilis-
 180 tic context-free grammar [6] provided by the researcher, that represents the
 181 possible choices for algorithms and hyperparameters in a machine learning
 182 pipeline. The grammar includes a probability distribution σ , which is ini-
 183 tially uniform for each of the decisions. A pipeline is obtained by sampling
 184 the grammar using this probability distribution. The optimisation process
 185 is based on probabilistic grammatical evolution [20], designed to incremen-
 186 tally adjust the probability distribution based on the marginal distribution
 187 of the best pipelines in each iteration. With a sufficiently large number of
 188 iterations, this process converges towards generating the best pipelines [13].

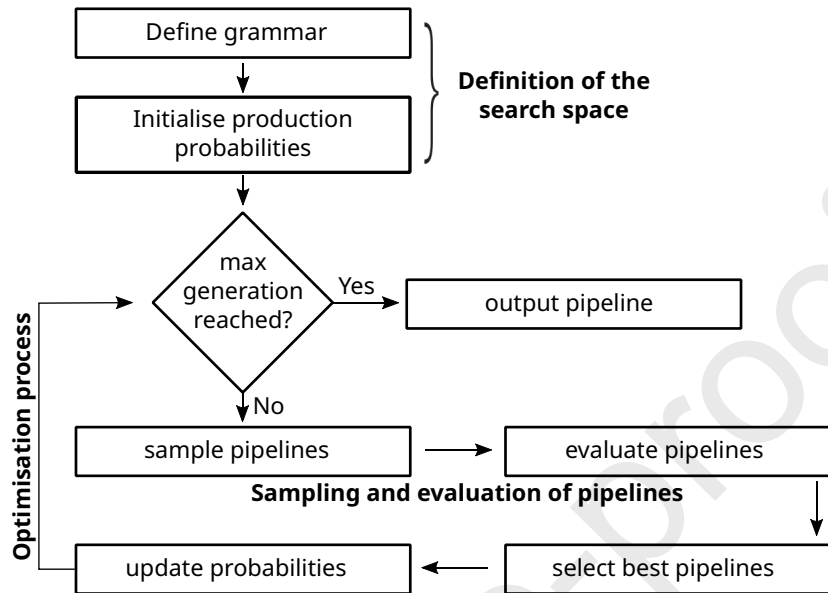


Figure 2: High-level representation of the Hierarchical Machine Learning optimisation process.

189 3.1. Definition of the search space

190 We define the search space for a specific machine learning problem as the
 191 set of all possible pipelines that solve the problem. Two different pipelines
 192 can use different algorithms for one specific step, or can use the same algo-
 193 rithm but with different values for some of its hyperparameters. Also, some
 194 steps can be combined under a higher-level step, forming a hierarchy of steps.
 195 For example, in NLP, a general “preprocessing” step can be composed of sub-
 196 steps such as “cleaning” and “dimensionality reduction”, which are in turn
 197 composed of more specific steps, until reaching algorithms and hyperparam-
 198 eters. To represent all possible pipelines for a given problem, and capture
 199 their hierarchical nature, we propose the use of context-free grammars. The
 200 grammar also defines the importance of each step of the pipelines. High-level
 201 decisions (e.g., using linear or non-linear classifiers) are represented closer to
 202 the start symbol of the grammar, while low-level decisions (e.g., the value of
 203 a regularisation factor) are closer to the end of the grammar. This way the
 204 researcher imposes a structure in the search space in which different pipelines
 205 can be compared.

206 As an illustrative example, the grammar in Figure 3 defines a small

- (1) $\langle \text{Pipeline} \rangle := \langle \text{Prep} \rangle \langle \text{Vect} \rangle \langle \text{Reduct} \rangle \langle \text{Class} \rangle$
- (2) $\langle \text{Prep} \rangle := \text{yes} \mid \text{no}$
- (3) $\langle \text{Vect} \rangle := \text{TF} (\langle \text{tf-ngram} \rangle) \mid \text{CV} (\langle \text{cv-ngram} \rangle)$
- (4) $\langle \text{tf-ngram} \rangle := i(1, 3)$
- (5) $\langle \text{cv-ngram} \rangle := i(1, 3)$
- (6) $\langle \text{Reduct} \rangle := \text{none} \mid \text{SVD}$
- (7) $\langle \text{Class} \rangle := \text{NB} \mid \text{LR} (\langle \text{reg} \rangle, \langle \text{penalty} \rangle) \mid \text{SVM} (\langle \text{kernel} \rangle)$
- (8) $\langle \text{reg} \rangle := f(0.01, 10)$
- (9) $\langle \text{penalty} \rangle := l1 \mid l2$
- (10) $\langle \text{kernel} \rangle := \text{linear} \mid \text{rbf}$

Figure 3: Grammar defining a set of pipelines for text classification. Each production represents a choice of the pipeline, which can be between different algorithms or their parameters.

207 search space to solve a text classification problem. The pipelines defined
 208 include: a text preprocessing phase —**Prep**—; different vectorial represen-
 209 tations —**Vect**—; an optional dimensionality reduction —**Reduct**—; and,
 210 different classifiers —**Class**—. The preprocessing step —production 2—
 211 consists of an optional application of stop-words removal. The vectorisation
 212 step —production 3— consists of a bag-of-word representation either using
 213 counts —**CV**— or TF-IDF weights —**TF**—. In each case, the grammar also
 214 specifies an n-gram size, an integer value between 1 and 3 —productions 4
 215 and 5—. Dimensionality reduction —production 6—, if applied, uses a stan-
 216 dard singular value decomposition (SVD). For classification —production
 217 7—, three algorithms are available, either naive bayes —**NB**—, logistic re-
 218 gression —**LR**— or support vector machines —**SVM**—. In the case of logistic
 219 regression —production 8—, two different regularisation functions are avail-
 220 able, either L_1 or L_2 , and in each case the regularisation factor is allowed to
 221 take values from 0.01 to 10 —production 9—. Finally, in the case of SVM,
 222 either a linear kernel or a radial basis function kernel (RBF) is used.

223 3.2. Sampling and evaluation of pipelines

224 The grammar defined has an attached probabilistic model σ , initialised
 225 as a collection of uniform distributions for each production of the grammar.
 226 Hence, every production has an assigned probability distribution π_i with
 227 parameters that change in time (as explained in Section 3.3). From this
 228 model σ a random pipeline is sampled using a top-down process. Starting at
 229 the root of the grammar (**Pipeline**), one of the available options is randomly

230 chosen according to the corresponding distribution π_0 . In the exemplary
 231 grammar shown in Figure 3, the initial production has only one option. Then,
 232 for each symbol in the selected production, the same sampling process is
 233 applied recursively. For example, in the case of production 7 —**Class**—,
 234 considering initial uniform probabilities, one of the three options (i.e., naive
 235 Bayes, logistic regression or SVM) is randomly selected. If naive Bayes —
 236 **NB**— is selected, the sampling process ends. However, if logistic regression —
 237 **LR**— is selected, then the sampling process proceeds recursively to production
 238 8. Figure 4 shows an example, illustrating the top-down sampling process
 239 and the final pipeline obtained for an arbitrary random selection.

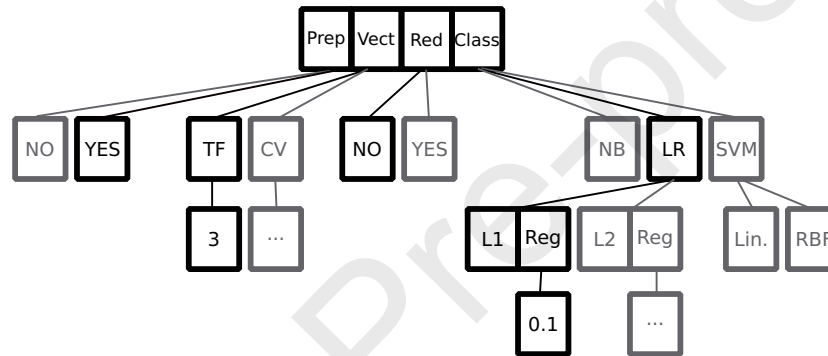


Figure 4: Sampling a pipeline from the exemplary grammar. The path that defines this specific pipeline is highlighted in **bold**.

240 During the top-down sampling process, each production of the gram-
 241 mar is interpreted with a different probabilistic distribution, according to
 242 its type. Productions can represent categorical decisions (i.e., one out of
 243 many), discrete or continuous values for hyper-parameters, or Boolean deci-
 244 sions (actually a special case of categorical decisions, but treated separately
 245 for convenience). According to the intended production type, the following
 246 distributions are used:

247 **Boolean values:** A Bernoulli distribution is used, with a parameter p that
 248 models the probability of obtaining **True**. A random value is generated
 249 by sampling from a random uniform value $u \in [0, 1]$, and the Boolean
 250 output is conditioned to $u \geq p$. In the previous example of text opinion
 251 mining, we can decide whether to delete stop-words.

252 **Categorical values:** Similar to the previous case, a multinomial Bernoulli

253 distribution is used, with parameters p_i , one for each option. When
 254 sampling, a random uniform value $u \in [0, 1]$ is generated and the p_i
 255 weights are accumulated until summing u , i.e., the result is the largest
 256 index k such that $\sum_{i=0}^k p_i \leq u$. This type of value is used for repre-
 257 senting a selection between a finite number of options where order is
 258 irrelevant, and only their relative weights are considered. For example,
 259 in the case of text opinion mining, we can decide between different
 260 classifiers: logistic regression, naive Bayes or support vector machines.

261 **Continuous values:** A normal distribution with mean μ and variance σ
 262 is used, plus two values a, b that correspond to the interval limits. When
 263 sampling, a Gaussian variable $N(\mu, \sigma)$ is generated (via the inverse
 264 method) and clamped to the limits a, b defined in the production. This
 265 type of value is used for representing continuous values. For exam-
 266 ple, in the case of text opinion mining we can decide the value of the
 267 regularisation parameter for logistic regression.

268 **Discrete values:** A normal distribution is also used as in the previous case,
 269 and the end results are rounded to the nearest integer. This distribu-
 270 tion encourages sampling to converge to a contiguous range of values
 271 as the search progresses, something that is not possible with a tradi-
 272 tional multinomial uniform distribution. This type of value is used for
 273 representing a selection between a finite number of options ordered by
 274 their numeric value. For example, in the case of text opinion mining,
 275 we can decide between different n-gram lengths: uni-grams, bi-grams
 276 or tri-grams.

277 Every pipeline sampled from the grammar is assigned a fitness value ac-
 278 cording to the machine learning problem at hand. For example, in the case of
 279 supervised learning, each pipeline presumably involves a classification or re-
 280 gression algorithm at some steps. Thus, each pipeline will be trained and
 281 evaluated on a suitable dataset, measuring precision, F_1 , R^2 , or another rele-
 282 vant metric. Even though statistically significant comparisons require several
 283 cross-validation iterations for each pipeline, a smaller number of evaluations
 284 per pipeline suffices during the optimisation process. This is due to the
 285 stochastic nature of the optimisation process explained in Section 3.3, which
 286 will tend to smooth away small random differences in pipeline fitness. In-
 287 spired by similar considerations in the design of ML-Plan [29], we recommend
 288 between 3 and 10 Monte Carlo cross-validation evaluations per pipeline.

289 3.3. Optimisation process

290 The optimisation process is based on a continuous optimisation technique
 291 with Grammatical Evolution for probabilistic context-free grammars by Kim
 292 and Ahn [20]. The process consists of a generation and evaluation loop using
 293 a grammar G suitable for the machine learning problem at hand, as defined
 294 in Section 3. In each iteration, a number of pipelines N is generated by
 295 sampling the grammar G according to the probabilities σ assigned to each
 296 production. These probabilities are initiated with a uniform distribution σ_0
 297 across all productions. Every pipeline is evaluated (which may consist of
 298 a single training/test execution, or a more complex cross-validation), and
 299 the best performing pipelines are used to modify σ , so as to maximise the
 300 probability of generating them.

301 A learning factor $\alpha \in [0, 1]$ is used to interpolate between the previ-
 302 ous distribution σ and the marginal distribution σ^* computed from the best
 303 pipelines. This learning factor controls the balance between the exploration
 304 and exploitation. A large α changes σ more rapidly, increasing exploitation,
 305 while a smaller value promotes more exploration. Recommended values for
 306 this parameter are 0.01 to 0.05. Algorithm 1 summarises the optimisation
 307 process for an arbitrary grammar G .

308 4. Results

309 This section presents two different evaluation scenarios to demonstrate
 310 the applicability and effectiveness of our proposal. The first scenario (Sec-
 311 tion 4.1) presents a comparison between HML-Opt and five AutoML frame-
 312 works in several standard machine learning problems. The purpose of this
 313 scenario is to demonstrate that HML-Opt can be used as an out-of-the-
 314 box AutoML tool that is competitive with the state-of-the-art. The second
 315 scenario (Section 4.2) applies HML-Opt to a complex knowledge discovery
 316 problem that includes not only selecting different algorithms, but also mod-
 317 elling high-level tasks that can be performed in different orders. The purpose
 318 of this scenario is to show the performance of HML-Opt in a real-life appli-
 319 cation, with a much larger search space, that includes NLP tasks, several
 320 shallow classifiers, neural networks, and different semantic features. Existing
 321 AutoML tools are not sufficiently expressive to model the space of possible
 322 pipelines required in this task. Hence, in this scenario, HML-Opt can only
 323 be compared with human-designed pipelines.

Algorithm 1 Grammatical Evolution for HML-Opt

$N \leftarrow$ size of population
 $n \leftarrow$ number of individuals selected each iteration
 $\alpha \leftarrow$ learning factor
 $G \leftarrow$ grammar for specific machine learning pipeline
 $\sigma_0 \leftarrow$ initial probabilities (uniform)
 $f \leftarrow$ fitness function (training and evaluation of a pipeline)
 Best $\leftarrow \square$

for each iteration i **do**
 $P_i \leftarrow$ sample population from G, σ_{i-1}
 for each solution $S \in P_i$ **do**
 $f(S) \leftarrow$ calculate fitness of S (evaluate pipeline)
 end for
 Best $\leftarrow \operatorname{argmax}_{S \in P \cup \{\text{Best}\}} \{f(S)\}$
 $P_i^* \leftarrow$ best n individuals of P_i
 $\sigma_i^* \leftarrow$ compute marginal distribution of P_i^*
 $\sigma_i \leftarrow \alpha \sigma_i^* + (1 - \alpha) \sigma_{i-1}$

end for
return Best

324 *4.1. Comparison with other AutoML frameworks*

325 One practical use of the AutoML paradigm is the design of out-of-the-box
 326 solutions that can be readily applied to a large number of similar problems.
 327 Using the `scikit-learn` library [33], we define a grammar that incorporates
 328 several different classifiers and preprocessing steps. This grammar is in part
 329 inspired by the design of Auto-Sklearn [12] and includes 4 data preprocessing
 330 steps, 7 different feature preprocessing techniques and 17 different classifi-
 331 cation algorithms. With this grammar, we implemented a `scikit-learn`
 332 compatible classifier that, in the same spirit of Auto-Sklearn, automatically
 333 chooses the best combination for a black-box classification problem. The
 334 grammar and complete source code to replicate the experimental setup for
 335 this example is available online². Figure 5 shows a graphical representation
 336 of the pipeline space defined by the grammar, not taking into consideration

²<http://bit.ly/hml-opt-sklearn>

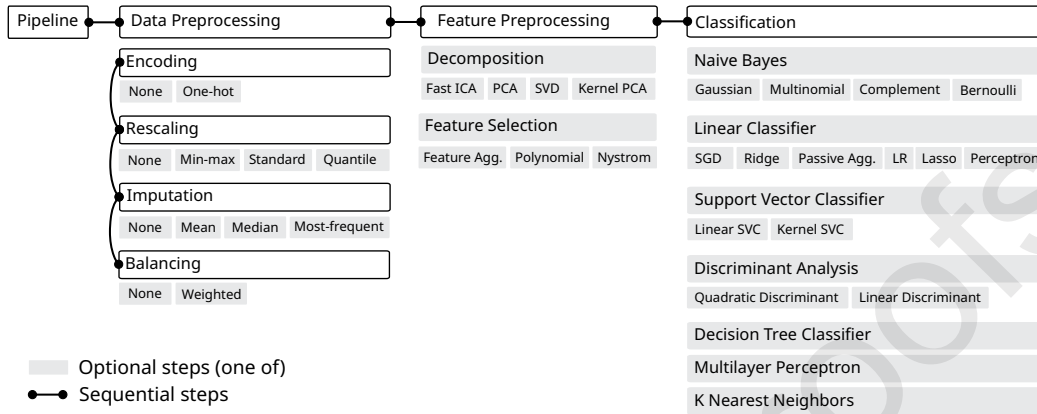


Figure 5: Graphical representation of the pipeline space defined by the *scikit-learn*-inspired grammar. Possible pipelines consist of three steps: data preprocessing, feature preprocessing and classification. Each step is recursively composed of sequential or optional sub-steps.

337 the hyperparameters of each technique.

338 Table 2 shows a comparison of our proposal, HML-Opt (Sklearn), to other
 339 similar AutoML frameworks, for different classic datasets. This experimen-
 340 tal setup is based on the results reported by Mohr et al. [29], and closely
 341 follows their proposed methodology. Specifically, we report the mean accu-
 342 racy across 20 independent executions of HML-Opt (Sklearn) in each dataset,
 343 with a random 70% of the data for optimisation and the remaining 30% as a
 344 hold-out test set, where the final accuracy is measured. Each execution has
 345 a timeout of one hour, and performs 5 Monte Carlo cross-validation evalu-
 346 ations (70%/30%) for each pipeline, with individual timeouts of 5 minutes
 347 per evaluation. In all cases the population size is 20 with a selection of the
 348 best 5 and a learning rate of 0.05.

349 The objective of this comparison is to show that our approach can be used
 350 in the same manner as other AutoML frameworks to design a high-performing
 351 black-box classifier, capable of obtaining competitive results. Three datasets
 352 from the UCI repository [25] are used in this evaluation: the *Cars* dataset [46]
 353 is a multivariate classification problem with 1,728 examples and 6 categorical
 354 features; the *Credit G.* dataset³ is also a multivariate classification problem
 355 with 1,000 examples and 20 categorical or integer features; and, the *Wine*

³[https://archive.ics.uci.edu/ml/datasets/statlog+\(german+credit+data\)](https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data))

Dataset	Cars	Credit G.	Wine Q.
ML-Plan (Weka)	98.73 (± 0.56)	74.46 (± 1.28)	67.38 (± 0.91)
Auto-WEKA	99.34 (± 0.38)	73.50 (± 2.32)	66.31 (± 1.90)
ML-Plan (Sklearn)	99.66 (± 0.51)	75.44 (± 2.53)	67.47 (± 1.31)
Auto-Sklearn-v	98.62 (± 0.67)	74.05 (± 1.89)	63.17 (± 1.27)
Auto-Sklearn-we	98.74 (± 0.53)	74.61 (± 0.88)	64.13 (± 1.18)
TPOT	99.63 (± 0.33)	76.09 (± 2.22)	67.06 (± 1.09)
HML-Opt (Sklearn)	99.92 (± 0.11)	75.03 (± 2.55)	<u>59.99</u> (± 2.74)

Table 2: Comparison of HML-Opt to other AutoML frameworks for three classic machine learning datasets in terms of mean and standard deviation of accuracy. Values for other frameworks were obtained from Mohr et al. [29]. Our results are highlighted in **bold** or underlined if they are statistically significantly **larger** or smaller than the rest, respectively, according to a Student’s t -test with a p -value < 0.05 .

356 Q . dataset [8] is a multivariate classification problem with 4,898 examples
 357 and 12 continuous features.

358 Using the same underlying grammar and optimisation parameters for all
 359 datasets, our approach obtains the best result in the Cars dataset, a statis-
 360 tically similar result in the Credit German dataset and a lower result in the
 361 Wine Quality dataset. In the latter dataset, we noticed that the one hour
 362 timeout was not sufficient for the evolutionary optimisation to achieve con-
 363 vergence, resulting in less than optimal performance. Since the search space
 364 and underlying algorithms in HML-Opt (Sklearn) are very similar to two
 365 of the comparisons that also use `scikit-learn`, it can be inferred that the
 366 differences in performance are due to differences either in the grammar defini-
 367 tion or in the exploration of the possible pipelines. These results demonstrate
 368 that HML-Opt can be used as an effective out-of-the-box AutoML tool in
 369 several machine learning scenarios.

370 4.2. Knowledge Extraction in Spanish eHealth text

371 In this section we present an example of the use of HML-Opt for solving
 372 a more complex machine learning problem as has been described by Estevez-
 373 Velarde et al. [11]. The problem selected is the eHealth Knowledge Discovery
 374 challenge presented in the TASS 2018 workshop [28]. The purpose of this
 375 case study is to show how HML-Opt can be extended beyond typical black-
 376 box machine learning problems and applied to complex domain problems,
 377 with the design of a suitable problem-specific grammar.

378 In the eHealth-KD challenge, participants were asked to submit a system
 379 that can identify relevant key phrases in a set of Spanish language docu-
 380 ments, classify them into two categories, and then detect semantic relations
 381 among them. The problem is divided into three different subtasks, that can
 382 be performed sequentially, or combined in different ways: (A) detecting enti-
 383 ties in text; (B) classifying them; and (C) detecting semantic relations among
 384 them. Figure 6 shows an example of sentences from the challenge, with their
 385 corresponding annotations. For this problem, a corpus with 1,173 sentences
 386 is provided [34] with a total of 13,113 examples, of which 3,310 are used for
 387 evaluation, and the rest for training. Results for this challenge show that a
 388 variety of techniques obtain competitive results, including: classic machine
 389 learning; deep learning; natural language processing; and, knowledge-based
 390 approaches. An effective exploration of the space of suitable pipelines re-
 391 quires mixing technologies from different software frameworks and knowledge
 392 areas. This makes the eHealth-KD challenge a complex scenario where exist-
 393 ing AutoML tools that provide out-of-the-box classifiers cannot be directly
 394 applied without considerable customisation.

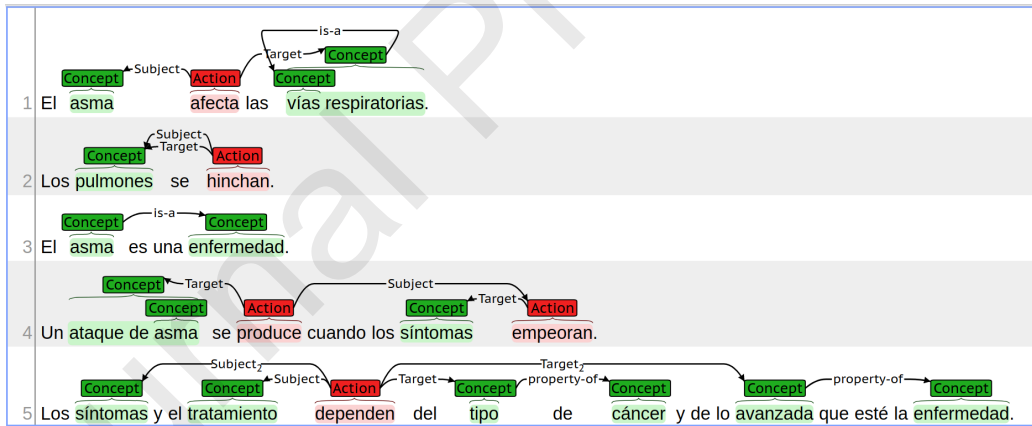


Figure 6: Example of annotated sentences in the eHealth-KD challenge. Adapted from Piad-Morffis et al. [34].

395 For this problem, we define a grammar that comprises several possible
 396 pipelines submitted to the official challenge evaluation. However, given the
 397 length of the grammar, it is not reported in full in this work although some
 398 insights on its characteristics are subsequently discussed. Our grammar con-
 399 siders all the 4 possible options for performing the three tasks sequentially

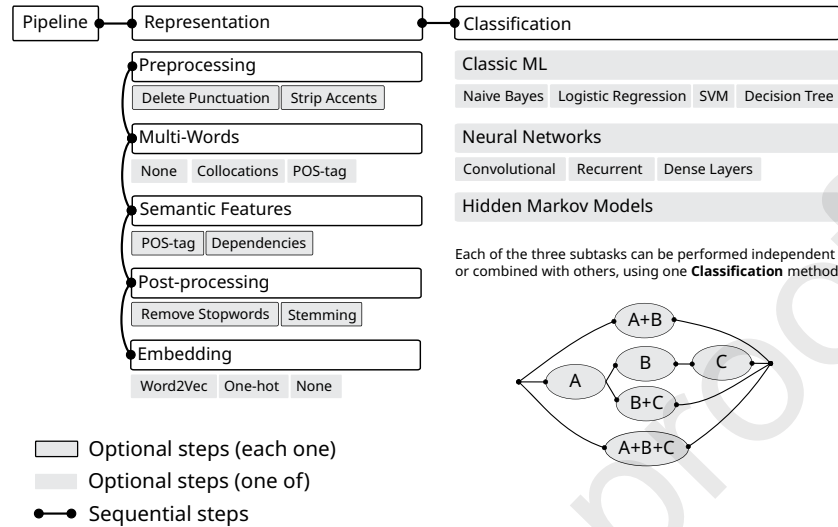


Figure 7: Graphical representation of the pipeline space defined by the eHealth-KD grammar. Possible pipelines consist of a representation step and a classification step where each subtask (i.e., A, B, and C) can be performed either independently or combined with each other.

400 combined. In terms of features, several options are defined, including syntac-
 401 tic features, word embedding techniques and knowledge-based features using
 402 the `nlTK` and `gensim` libraries. For the classification part we include both
 403 classic machine learning techniques from `scikit-learn` as well as deep learn-
 404 ing architectures from `keras` [7]. The grammar and complete source code for
 405 this case study is available online⁴. Figure 7 shows a graphical representation
 406 of the pipeline space defined by this grammar. Due to the hierarchical nature
 407 of the grammar, not all details can be captured in this figure.

408 The optimisation process was performed for a total of 60 generations,
 409 with a population size of 50 pipelines per generation, a selection factor of
 410 the best 20% and a learning factor of 0.05. The average time required to
 411 evaluate each pipeline was 5.17 minutes for a total of approximately 257
 412 hours of computation in a single computational node of commodity hardware.
 413 Long-running pipelines were stopped after a 10 minute timeout. The best
 414 performing pipeline was encountered at the eighteenth generation.

415 Table 3 shows the best results obtained by HML-Opt for this problem.

⁴<http://bit.ly/hml-opt-ehealthkd>

416 For comparative purposes, the results obtained by other authors are taken
 417 from the TASS 2018 Overview [28]. We only evaluate HML-Opt on the
 418 proposed evaluation scenario that comprises all three subtasks. Hence, only
 419 authors that participate in this scenario are included. As Table 3 shows, in
 420 this problem HML-Opt obtains results competitive with the state-of-the-art.

Techniques	System	F_1
AutoML (NLP+LR+DT+SVM)	HML-Opt (TASS)	75.4
Bi-LSTM + CRF	Extended NeuroNER [36]	74.4
Knowledge-based	SINAI [27]	71.0
NLP + CRF + LR	UPF-UPC [32]	68.1
Convolutional NN	LABDA [38]	31.0

Table 3: Comparison of HML-Opt to other techniques used on the eHealth-KD dataset, specifically in the evaluation scenario No. 1. Following the dataset evaluation methodology, a single evaluation on the official test set was performed.

421 Interestingly, the final pipeline obtained by HML-Opt in this example is
 422 composed of classic algorithms: SVM, logistic regression and decision trees.
 423 In terms of representations, HML-Opt selected one pipeline that included
 424 one-hot encoding with bi-grams, stemming, and POS-tagging as an addi-
 425 tional syntactic feature. By contrast, other authors rely on deep learning
 426 architectures, word and character embedding, and external knowledge bases.
 427 The simple but effective pipeline discovered by HML-Opt shows that an ad-
 428 equate optimisation and automatic parameter tuning can out-perform more
 429 complex but insufficiently tuned strategies that are hand-designed.

430 It is important to emphasise that in this case study it is impossible to
 431 compare with other AutoML tools since, to the extent of our knowledge, none
 432 of the existing technologies allow for the degree of versatility necessary in this
 433 NLP problem. To solve the aforementioned problem, we need to combine
 434 technologies from multiple machine learning frameworks, e.g., embeddings
 435 from `gensim`, structured features from `sklearn` and learning architectures
 436 from `keras`. HML-Opt leverages expert domain knowledge expressed in the
 437 form of a problem-specific grammar, and is thus able to combine a variety of
 438 tools and frameworks in a non-trivial manner to produce pipelines that are
 439 specifically designed for this problem.

440 5. Discussion

441 In this section, we present a high-level analysis of the experimental re-
442 sults and discuss key characteristics of HML-Opt. Unlike other proposals,
443 our technique is not restricted to particular types of pipelines, such as neu-
444 ral networks, or shallow classifiers. Since the grammar defines the space
445 of possible experimentation, anything can be included, such as natural lan-
446 guage preprocessing techniques or knowledge bases. Our proposal optimises
447 at many different levels, from high-level decisions about classes of algorithms
448 to low-level decisions such as hyperparameters (e.g., regularisation factors).

449 The fundamental basis of our proposal is that researchers can define
450 a problem-specific grammar that encodes their domain knowledge about a
451 problem. This way, researchers define which pipelines are interesting to ex-
452 plore, while HML-Opt evaluates and searches for the best pipeline. However,
453 for classic problems, or in a real life scenario, having access to an out-of-the-
454 box classifier can be desirable. In these cases, our approach can be used
455 similar to other AutoML frameworks, to predefine a set of “black-box” solu-
456 tions, i.e., problem-specific grammars that are tailored for specific problems.
457 Researchers or practitioners can start by using these black-box solutions,
458 but still have the option to take a look inside the black-box and modify the
459 grammar for their specific purposes, where necessary.

460 The next step in the automation of machine learning would be to au-
461 tomate the process of defining a suitable grammar for a specific problem.
462 In our proposal, researchers have to manually design the appropriate gram-
463 mar based on their knowledge and previous experience about the problem.
464 However, analysing the case studies presented in Section 4, we observed that
465 although the two grammars are different, they also share similarities. For
466 example, they consist of an initial production that defines high-level steps,
467 which in turn are defined into simpler steps, and so on, down to algorithms
468 and hyperparameter values. Furthermore, they all define some sort of pre-
469 processing and representation steps, even if they are named differently and
470 composed of different sub-steps, and a final classification step. These and
471 other regularities could lay the foundation for building a generic knowledge
472 base of machine learning techniques that would help automate the definition
473 of problem-specific grammars.

474 Approaching the AutoML problem from the perspective of its underlying
475 optimisation model, two interesting considerations become apparent: the
476 hierarchical nature of the solutions and the multi-objective nature of the

477 performance metric. Our proposal explicitly represents solutions as a hierar-
 478 chical decision process, in the same way as other some alternative AutoML
 479 approaches (i.e., Recipe and ML-Plan). This allows modelling and account-
 480 ing for the fact that some decisions have a greater influence on the pipeline
 481 optimisation, for instance, selecting between linear, non-linear or tree-based
 482 classifiers supersedes selecting specific hyperparameter values. We argue that
 483 as machine learning pipelines grow in complexity and involve algorithms from
 484 different libraries and technologies, a hierarchical conceptualisation will be
 485 even more relevant. With respect to the performance metric, most current
 486 AutoML approaches focus on a single objective function, e.g., accuracy, pre-
 487 cision, F_1 , in correspondence with the machine learning problem at hand.
 488 However, in practical scenarios, it may be necessary to balance different
 489 performance metrics, including also time and memory usage, and more sub-
 490 jective qualities such as explainability of the models. The simplest approach
 491 (taken by HML-Opt and others) of optimising a main metric subject to re-
 492 strictions of time and memory is insufficient in a scenario where the end-user
 493 has to decide about practical issues such as deploying these pipelines in a
 494 production system. TPOT considers this problem from the multi-objective
 495 approach by jointly optimising both accuracy and model complexity (in terms
 496 of pipeline length). However, future approaches will need to consider this is-
 497 sue in-depth, potentially including the end-user in the pipeline evaluation
 498 loop.

499 Analysing the evolution of the performance of pipelines provides some
 500 insight on the complexity of a given machine learning problem. As an il-
 501 lustrative example, Figure 8 (left) shows the evolution of the average and
 502 optimal accuracy in the case study presented in Section 4.2. The average
 503 accuracy steadily increased, evidencing the optimisation process. The rela-
 504 tively low average accuracy at the beginning is in large part influenced by the
 505 existence of *invalid pipelines*. Invalid pipelines are encountered in two cases:
 506 when a pipeline exceeds a timeout predefined by the researcher, or when
 507 unpredictable run-time errors occur, such as out-of-memory errors provoked
 508 by an unfeasible combination of hyperparameters. These circumstances are
 509 often impossible to predict by the researcher beforehand and as such cannot
 510 be accounted for in the grammar. Invalid pipelines are scored “0” for fit-
 511 ness, to steer the optimisation process away from them. As Figure 8 (right)
 512 shows, the number of invalid pipelines decreases steadily, which is an indica-
 513 tion that the probabilistic model is penalising the characteristics that appear
 514 consistently in invalid pipelines.

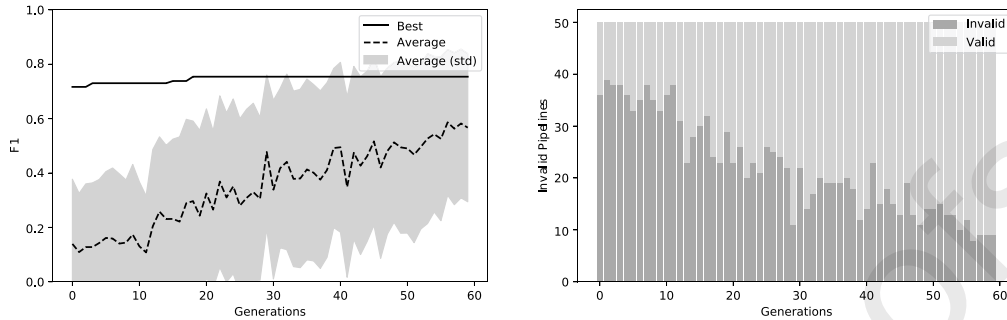


Figure 8: **Left:** Evolution of the average and optimal accuracy of pipelines in each iteration. **Right:** Proportion of invalid vs. valid pipelines found in each iteration. Results are from the case study presented in Section 4.2.

515 During the optimisation process a significant amount of data is gener-
 516 ated in each iteration about the best pipelines and their characteristics. In
 517 principle, this information is used to update the probabilistic model so that
 518 new sampled pipelines are similar to the best pipelines already evaluated.
 519 However, this information can provide additional insights when aggregated
 520 for the whole optimisation process, by analysing the characteristics (i.e., the
 521 productions of the grammar) that consistently appear in the best pipelines.

522 Specifically for the case study presented in section 4.2, we analysed the
 523 classification algorithms used in each subtask by training a linear regression
 524 on all the valid pipelines evaluated during the optimisation process to predict
 525 the accuracy of a pipeline given its characteristics. The features used in the
 526 linear regression are obtained directly from each pipeline’s derivation tree by
 527 considering all the decisions involved in a pipeline. Every path from the root
 528 to a leaf in the derivation tree becomes a feature. As an example, if one
 529 specific pipeline applies Logistic Regression with a regularisation factor of
 530 0.3 for subtask A, then a feature `A.Classifier.LR.regularisation=0.3` is
 531 generated. Our grammar involves a total of 242 such features, which account
 532 for all possible decisions for each pipeline, from high-level task ordering to
 533 hyperparameter values.

534 The coefficients assigned to each characteristic indicate which character-
 535 istic has a positive or negative impact on the pipeline’s fitness. Rather
 536 than exclusively analysing the best pipeline, these coefficients indicate gen-
 537 eral trends about which algorithms are consistently more effective in the
 538 majority of pipelines. For instance, Table 4 shows these coefficients only

539 for the classification algorithms used in each subtask, aggregated across all
 540 hyperparameter combinations. In this case, the best pipeline uses Logistic
 541 Regression, SVM and Decision Tree as classifiers (see Section 4.2). However,
 542 Table 4 shows that Neural Networks, on average, tend to outperform other
 543 algorithms, even though no specific pipeline using Neural Networks was the
 544 optimum.

Algorithm	Coefficients by subtask		
	A	B	C
Logistic Regression	0.0014	0.0074	-0.0127
Neural Networks	0.0361	0.0320	-0.0009
Support Vector Machines	-0.0157	-0.0221	0.0018
Decision Tree	0.0165	0.0223	0.0024
Naive Bayes	0.0222	0.0128	-
Markov Models	-0.0081	-	-

Table 4: Relevance of the classification algorithms in each subtask (A, B and C) in the case study presented in Section 4.2. The highest coefficient in each subtask is highlighted. Missing values correspond to combinations that were not evaluated in any valid pipeline.

545 A similar analysis can shed light on the relevance of the NLP interme-
 546 diate tasks. Table 5 shows a subset of the features corresponding to the
 547 representation phase and their corresponding coefficients. The most relevant
 548 tasks appear to be the computation of POS-tag and dependency labels, while
 549 the use of embeddings, counter-intuitively seems to provide a marginal neg-
 550 ative effect. Interestingly, the optimal solution found by HML-Opt does not
 551 necessarily comprise the components with the highest weights. For example,
 552 the optimal solution applies the `postag` strategy for multi-word represen-
 553 tation and performs stemming but does not perform stop-words removal,
 554 directly contradicting what these coefficients suggest. Researchers can use
 555 these insights to better understand the characteristics of the problem and
 556 design more effective pipelines. However, caution must be exercised, since
 557 this analysis rests on the assumption that different decisions in a machine
 558 learning pipeline have an independent impact on the overall performance of
 559 the model, which can be characterised by a simple linear relationship.

560 The analysis presented in this section is only illustrative of the type of
 561 insights that can be obtained when both a probabilistic model and a hi-
 562 erarchical representation of a machine learning pipeline are available. This

NLP Subtask	Strategy	Coefficient
Embedding	none	0.015
Embedding	onehot	-0.030
Embedding	wordVec	0.006
MultiWords	collocations	-0.032
MultiWords	none	0.014
MultiWords	postag	0.009
PostProcessing.Stemming	no	-0.003
PostProcessing.Stemming	yes	-0.005
PostProcessing.StopWords	no	-0.010
PostProcessing.StopWords	yes	0.001
Preprocessing.DelPunt	no	-0.001
Preprocessing.DelPunt	yes	-0.008
Preprocessing.StripAccents	no	0.001
Preprocessing.StripAccents	yes	-0.010
SemanticFeatures.Dependencies	no	-0.013
SemanticFeatures.Dependencies	yes	0.003
SemanticFeatures.PosTag	no	-0.025
SemanticFeatures.PosTag	yes	0.016

Table 5: Overall relevance of the representation features for the case study presented in Section 4.2. The highest coefficient for each feature is highlighted.

563 information can also be used to design more intelligent search procedures that
564 employ past experience in similar problems. For example, the probabilistic
565 model can be initialised with a non-uniform distribution proportional to the
566 weights obtained in previous experimentation. Furthermore, by designing
567 grammars with productions that represent high-level characteristics of the
568 pipelines (e.g., separating linear and non-linear classifiers), this analysis can
569 also provide high-level insights on which characteristics are more important
570 in specific problem domains. An in-depth analysis of this experimental setup
571 has already been published [11].

572 In terms of performance, the cost of executing HML-Opt is dominated by
573 the cost of training and evaluating each individual pipeline. Since pipelines

574 for machine learning problems are often computationally intensive, the ad-
575 ditional cost of updating the probabilistic model and the bookkeeping as-
576 sociated with the optimisation process is negligible. Several strategies can
577 help reduce the computational cost. For instance, since each generation
578 performs N independent pipeline evaluations, HML-Opt can benefit from
579 the implementation of parallel and distributed computing strategies. An-
580 other cost saving strategy consists of monitoring some convergence metrics
581 of the underlying probabilistic model to provide an early stopping mecha-
582 nism. However, we believe that a meta-learning strategy that learns from
583 previous experience to seed the probabilistic model can provide even higher
584 cost saving benefits by comparing the new dataset that requires solving to
585 previously solved problems. Such a strategy would leverage every past exe-
586 cution of HML-Opt, potentially even by other researchers, as every pipeline
587 evaluation ever made will work to more efficiently solve all future experi-
588 ments.

589 6. Conclusions

590 This paper presents Hierarchical Machine Learning Optimisation (HML-
591 Opt), an AutoML framework for automatically finding a close to optimal
592 pipeline in a specific machine learning problem. This novel technique allows
593 researchers to evaluate a much higher number of experimental setups than
594 what is manually possible, given specific time frames and computational
595 resources. Moreover, a key and innovative design feature of HML-Opt is
596 to provide a declarative and expressive framework whereby a researcher can
597 define the desired space of possible pipelines and explore it effectively and
598 efficiently. This contrasts with existing alternatives that tend to force the
599 researcher into a predefined set of algorithms or possible architectures for
600 pipelines. Experimental evaluations indicate that our proposal is competitive
601 with other approaches for AutoML, while at the same time applicable to
602 novel scenarios that require integrating tools from different machine learning
603 frameworks.

604 AutoML techniques like the one presented in this research are a staple in
605 the current trend of democratising machine learning technologies. The tools
606 continually evolve to provide more powerful ready-to-use machine learning
607 pipelines. However, AutoML systems should not only attempt to replace
608 human experts, but rather to serve as complementary tools that enable re-
609 searchers and practitioners to quickly obtain better baselines and insights

610 on the most promising strategies. Armed with AutoML, we envision a fu-
611 ture where people and computers work together in solving the most pressing
612 problems of mankind.

613 Acknowledgements

614 **Funding:** This research has been supported by a Carolina Foundation
615 grant in agreement with University of Alicante and University of Havana.
616 Moreover, it has also been partially funded by both aforementioned uni-
617 versities, the Generalitat Valenciana (*Conselleria d'Educació, Investigació,*
618 *Cultura i Esport*) and the Spanish Government through the projects LIVING-
619 LANG (RTI2018-094653-B-C22) and SIIA (PROMETEO/2018/089, PROM-
620 ETEU/2018/089).

621 This manuscript has been greatly improved by the valuable comments of
622 anonymous reviewers.

- 623 [1] F. Assuncao, N. Lourenco, P. Machado, and B. Ribeiro. Automatic
624 generation of neural networks with structured Grammatical Evolution.
625 In *2017 IEEE Congr. Evol. Comput. CEC 2017 - Proc.*, pages 1557–1564.
626 IEEE, 2017. ISBN 9781509046010. doi: 10.1109/CEC.2017.7969488.
- 627 [2] L. Bottou. Large-scale machine learning with stochastic gradient de-
628 scent. In *Proc. COMPSTAT'2010*, pages 177–186. Springer, 2010.
- 629 [3] E.K. Burke, M. Gendreau, M. Hyde, G. Kendall, G. Ochoa, E. Özcan,
630 and R. Qu. Hyper-heuristics: A survey of the state of the art. *J. Oper.*
631 *Res. Soc.*, 64(12):1695–1724, 2013. ISSN 14769360. doi: 10.1057/jors.
632 2013.71. URL <https://doi.org/10.1057/jors.2013.71>.
- 633 [4] F. Caraffini, F. Neri, and M. Epitropakis. HyperSPAM: A study
634 on hyper-heuristic coordination strategies in the continuous domain.
635 *Inf. Sci. (Ny)*, 477:186–202, 2019. ISSN 00200255. doi: 10.1016/
636 j.ins.2018.10.033. URL [http://www.sciencedirect.com/science/
637 article/pii/S002002551830851X](http://www.sciencedirect.com/science/article/pii/S002002551830851X).
- 638 [5] A.R. Carvalho, F.M. Ramos, and A.A. Chaves. Metaheuristics for the
639 feedforward artificial neural network (ANN) architecture optimization
640 problem. *Neural Comput. Appl.*, 20(8):1273–1284, 2011. ISSN 09410643.
641 doi: 10.1007/s00521-010-0504-3.

- 642 [6] Z. Chi. Statistical properties of probabilistic context-free grammars.
643 *Comput. Linguist.*, 25(1):131–160, 1999. ISSN 08912017.
- 644 [7] F. Chollet. Keras. <https://keras.io>, 2015.
- 645 [8] P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis. Modeling
646 wine preferences by data mining from physicochemical properties. *Decis.*
647 *Support Syst.*, 47(4):547–553, 2009. ISSN 01679236. doi: 10.1016/j.dss.
648 2009.05.016.
- 649 [9] A.G.C. de Sá, W.J.G.S. Pinto, L.O.V.B. Oliveira, and G.L. Pappa.
650 RECIPE: A grammar-based framework for automatically evolving clas-
651 sification pipelines. In *Lect. Notes Comput. Sci. (including Subser.*
652 *Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, volume 10196
653 LNCS, pages 246–261. Springer, 2017. ISBN 9783319556956. doi:
654 10.1007/978-3-319-55696-3_16.
- 655 [10] A.M. De Silva and P.H.W. Leong. Grammatical evolution. *Springer-*
656 *Briefs Appl. Sci. Technol.*, 5(9789812874108):25–33, 2015. ISSN
657 21915318. doi: 10.1007/978-981-287-411-5_3.
- 658 [11] S. Estevez-Velarde, Y. Gutiérrez, A. Montoyo, and Y. Almeida-Cruz.
659 AutoML strategy based on grammatical evolution: A case study about
660 knowledge discovery from text. In *ACL 2019 - 57th Annu. Meet. Assoc.*
661 *Comput. Linguist. Proc. Conf.*, pages 4356–4365, Florence, Italy, 2020.
662 Association for Computational Linguistics. ISBN 9781950737482. URL
663 <https://www.aclweb.org/anthology/P19-1428>.
- 664 [12] M. Feurer, A. Klein, K. Eggenberger, J.T. Springenberg, M. Blum, and
665 F. Hutter. Efficient and robust automated machine learning. In *Adv.*
666 *Neural Inf. Process. Syst.*, volume 2015-Janua, pages 2962–2970, 2015.
- 667 [13] M. Hauschild and M. Pelikan. An introduction and survey of estimation
668 of distribution algorithms. *Swarm Evol. Comput.*, 1(3):111–128, 2011.
669 ISSN 22106502. doi: 10.1016/j.swevo.2011.08.003.
- 670 [14] F. Hutter, L. Kotthoff, and J. Vanschoren, editors. *Automated Machine*
671 *Learning: Methods, Systems, Challenges*. Springer, 2018.
- 672 [15] H. Jin, Q. Song, and X. Hu. Auto-Keras: Efficient Neural Architecture
673 Search with Network Morphism, 2018. ISSN 00086363 (ISSN).

- 674 [16] H. Josiński, D. Kostrzewa, A. Michalczyk, and A. Świtoński. The ex-
675 panded invasive weed optimization metaheuristic for solving continuous
676 and discrete optimization problems. *Sci. World J.*, 2014, 2014. ISSN
677 1537744X. doi: 10.1155/2014/831691.
- 678 [17] J. Józefowska and J. Weglarz. On a methodology for discrete-continuous
679 scheduling. *Eur. J. Oper. Res.*, 107(2):338–353, 1998. ISSN 03772217.
680 doi: 10.1016/S0377-2217(97)00346-9.
- 681 [18] P. Kerschke and H. Trautmann. Automated algorithm selection on con-
682 tinuous black-box problems by combining exploratory landscape anal-
683 ysis and machine learning. *Evol. Comput.*, 27(1):99–127, 2018. ISSN
684 15309304. doi: 10.1162/evco_a_00236. URL [https://doi.org/10.](https://doi.org/10.1162/evco_a_00236)
685 [1162/evco_a_00236](https://doi.org/10.1162/evco_a_00236).
- 686 [19] K. Khan and A. Sahai. A Comparison of BA, GA, PSO, BP and LM
687 for Training Feed forward Neural Networks in e-Learning Context. *Int.*
688 *J. Intell. Syst. Appl.*, 4(7):23–29, 2012. ISSN 2074904X. doi: 10.5815/
689 ijisa.2012.07.03.
- 690 [20] H.T. Kim and C.W. Ahn. A New Grammatical Evolution Based
691 on Probabilistic Context-free Grammar. In *Proc. 18th Asia Pacific*
692 *Symp. Intell. Evol. Syst. 2*, pages 1–12. Springer, 2015. doi: 10.1007/
693 978-3-319-13356-0_1.
- 694 [21] D.P. Kingma and J.L. Ba. Adam: A method for stochastic optimization.
695 *3rd Int. Conf. Learn. Represent. ICLR 2015 - Conf. Track Proc.*, 2015.
- 696 [22] B. Komer, J. Bergstra, and C. Eliasmith. Hyperopt-Sklearn: Auto-
697 matic Hyperparameter Configuration for Scikit-Learn. In *Proc. 13th*
698 *Python Sci. Conf.*, pages 32–37. Citeseer, 2014. doi: 10.25080/
699 majora-14bd3278-006.
- 700 [23] K. Li and J. Malik. Learning to optimize. *5th Int. Conf. Learn. Repre-*
701 *sent. ICLR 2017 - Conf. Track Proc.*, 2019.
- 702 [24] Z. Li, X. Xiong, Z. Ren, N. Zhang, X. Wang, and T. Yang. An Aggressive
703 Genetic Programming Approach for Searching Neural Network Struc-
704 ture Under Computational Constraints. *arXiv Prepr. arXiv1806.00851*,
705 2018. URL <http://arxiv.org/abs/1806.00851>.

- 706 [25] M. Lichman. UCI Machine Learning Repository
707 [http://archive.ics.uci.edu/ml]., 2013. URL [http://archive.ics.
708 uci.edu/ml](http://archive.ics.uci.edu/ml).
- 709 [26] H. Liu, K. Simonyan, O. Vinyals, C. Fernando, and K. Kavukcuoglu.
710 Hierarchical representations for efficient architecture search. *6th Int.
711 Conf. Learn. Represent. ICLR 2018 - Conf. Track Proc.*, 2018.
- 712 [27] P. López-Úbeda, M.C. Díaz-Galiano, M.T. Martín-Valdivia, and L.A.
713 Ureña-López. SINAI in TASS 2018 Task 3. Classifying actions and
714 concepts with UMLS on MedLine. In *CEUR Workshop Proc.*, volume
715 2172, pages 77–82, 2018.
- 716 [28] E. Martínez-Cámara, Y. Almeida-Cruz, M.C. Díaz-Galiano, S. Estévez-
717 Velarde, M. García-Cumbreras, M. García-Vega, Y. Gutiérrez,
718 A. Montejo-Ráez, A. Montoyo, R. Muñoz, A. Piad-Morffis, and
719 J. Villena-Román. Overview of TASS 2018: Opinions, health and emo-
720 tions. In *CEUR Workshop Proc.*, volume 2172 of *CEUR Workshop Pro-
721 ceedings*, pages 13–27, Sevilla, Spain, 2018. CEUR-WS.
- 722 [29] F. Mohr, M. Wever, and E. Hüllermeier. ML-Plan: Automated machine
723 learning via hierarchical planning. *Mach. Learn.*, 107(8-10):1495–1515,
724 sep 2018. ISSN 15730565. doi: 10.1007/s10994-018-5735-z. URL [https:
725 //doi.org/10.1007/s10994-018-5735-z](https://doi.org/10.1007/s10994-018-5735-z).
- 726 [30] R. Negrinho and G. Gordon. DeepArchitect: Automatically Designing
727 and Training Deep Architectures. *arXiv Prepr. arXiv1704.08792*, 2017.
728 URL <http://arxiv.org/abs/1704.08792>.
- 729 [31] R.S. Olson and J.H. Moore. TPOT: A Tree-Based Pipeline Optimization
730 Tool for Automating Machine Learning. In *Work. Autom. Mach. Learn.*,
731 pages 151–160, 2019. doi: 10.1007/978-3-030-05318-5_8.
- 732 [32] J.V. Palatresi and H.R. Hontoria. TASS2018: Medical knowledge dis-
733 covery by combining terminology extraction techniques with machine
734 learning classification. In *CEUR Workshop Proc.*, volume 2172, pages
735 89–95. CEUR-WS. org, 2018.
- 736 [33] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion,
737 O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Van-

- 738 derplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duch-
 739 esnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
 740 *Learning Research*, 12:2825–2830, 2011.
- 741 [34] A. Piad-Morffis, Y. Gutiérrez, S. Estévez-Velarde, Y. Almeida-Cruz,
 742 A. Montoyo, and R. Muñoz. Analysis of eHealth knowledge discovery
 743 systems in the TASS 2018 Workshop. *Proces. Leng. Nat.*, 62(0):13–20,
 744 2019. ISSN 19897553. doi: 10.26342/2019-62-1. URL <http://journal.sepln.org/sepln/ojs/ojs/index.php/pln/article/view/5947>.
 745
- 746 [35] E. Real, S. Moore, A. Selle, S. Saxena, Y.L. Suematsu, J. Tan, Q.V.
 747 Le, and A. Kurakin. Large-scale evolution of image classifiers. In *34th*
 748 *Int. Conf. Mach. Learn. ICML 2017*, volume 6, pages 4429–4446, 2017.
 749 ISBN 9781510855144.
- 750 [36] R.M. Rivera-Zavala, P. Martínez, and I. Segura-Bedmar. A Hybrid Bi-
 751 LSTM-CRF model for knowledge recognition from ehealth documents.
 752 In *CEUR Workshop Proc.*, volume 2172, pages 65–70, 2018.
- 753 [37] K.O. Stanley and R. Miikkulainen. Evolving neural networks
 754 through augmenting topologies. *Evol. Comput.*, 10(2):99–127, jun
 755 2002. ISSN 10636560. doi: 10.1162/106365602320169811. URL
 756 [http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.](http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.28.5457&rep=rep1&type=pdf)
 757 [28.5457&rep=rep1&type=pdfhttp://www.mitpressjournals.org/](http://www.mitpressjournals.org/doi/10.1162/106365602320169811)
 758 [doi/10.1162/106365602320169811](http://www.mitpressjournals.org/doi/10.1162/106365602320169811).
- 759 [38] V. Suárez-Paniagua, I. Segura-Bedmar, and P. Martínez. LABDA at
 760 TASS-2018 Task 3: Convolutional neural networks for relation classifi-
 761 cation in Spanish eHealth documents. In *CEUR Workshop Proc.*, volume
 762 2172, pages 71–76, 2018.
- 763 [39] C. Thornton, F. Hutter, H.H. Hoos, and K. Leyton-Brown. Auto-
 764 WEKA: Combined selection and hyperparameter optimization of clas-
 765 sification algorithms. In *Proc. ACM SIGKDD Int. Conf. Knowl. Discov.*
 766 *Data Min.*, volume Part F1288, pages 847–855. ACM, 2013. ISBN
 767 9781450321747. doi: 10.1145/2487575.2487629.
- 768 [40] T. Tieleman and G. Hinton. Lecture 6.5-rmsprop: Divide the gradient by
 769 a running average of its recent magnitude. *COURSERA Neural networks*
 770 *Mach. Learn.*, 4(2):26–31, 2012.

- 771 [41] I. Tsoulos, D. Gavrilis, and E. Glavas. Neural network construction and
772 training using grammatical evolution. *Neurocomputing*, 72(1-3):269–277,
773 2008. ISSN 09252312. doi: 10.1016/j.neucom.2008.01.017.
- 774 [42] G. Waligóra. Discrete-continuous project scheduling with discounted
775 cash flows-A tabu search approach. *Comput. Oper. Res.*, 35(7):2141–
776 2153, 2008. ISSN 03050548. doi: 10.1016/j.cor.2006.09.022.
- 777 [43] B. Zoph and Q.V. Le. Neural architecture search with reinforcement
778 learning. *5th Int. Conf. Learn. Represent. ICLR 2017 - Conf. Track*
779 *Proc.*, 2019.
- 780 [44] B. Zoph, V. Vasudevan, J. Shlens, and Q.V. Le. Learning Transferable
781 Architectures for Scalable Image Recognition. In *Proc. IEEE Com-*
782 *put. Soc. Conf. Comput. Vis. Pattern Recognit.*, pages 8697–8710, 2018.
783 ISBN 9781538664209. doi: 10.1109/CVPR.2018.00907.
- 784 [45] B. Zoph, V. Vasudevan, J. Shlens, and Q.V. Le. Learning Transferable
785 Architectures for Scalable Image Recognition. *Proc. IEEE Comput.*
786 *Soc. Conf. Comput. Vis. Pattern Recognit.*, 2(6):8697–8710, 2018. ISSN
787 10636919. doi: 10.1109/CVPR.2018.00907.
- 788 [46] B. Zupan, M. Bohanec, I. Bratko, and J. Demsar. Machine learning by
789 function decomposition. In *ICML*, pages 421–429, 1997.