

Article

Optimized Distributed Hyperparameter Search and Simulation for Lung Texture Classification in CT Using Hadoop

Roger Schaer^{1,*}, Henning Müller^{1,2,†} and Adrien Depeursinge^{1,3,†}

¹ Information Systems Institute, University of Applied Sciences Western Switzerland (HES–SO), Techno–Pôle 3, 3960 Sierre, Switzerland; E–mail: henning.mueller@hevs.ch,

² University Hospitals and University of Geneva, Rue Gabrielle–Perret–Gentil 4, 1205 Geneva, Switzerland,

³ Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne 1015, Switzerland; E–mail: adrien.depeursinge@epfl.ch.

† These authors contributed equally to this work.

* Author to whom correspondence should be addressed; E–mail: roger.schaer@hevs.ch, Tel.: +41–27–606–9035

Version May 6, 2016 submitted to *Jimaging*. Typeset by *LaTeX* using class file *mdpi.cls*

1 **Abstract:** Many medical image analysis tasks require complex learning strategies to reach
2 a quality of image–based decision support that is sufficient in clinical practice. The analysis
3 of medical texture in tomographic images, for example of lung tissue, is no exception. Via a
4 learning framework, very good classification accuracy can be obtained but several parameters
5 need to be optimized. This article describes a practical framework for efficient distributed
6 parameter optimization. The proposed solutions are applicable for many research groups
7 with heterogeneous computing infrastructures and for various machine learning algorithms.
8 These infrastructures can easily be connected via distributed computation frameworks. We
9 use the Hadoop framework to run and distribute both grid and random search strategies for
10 hyperparameter optimization and cross–validations on a cluster of 21 nodes composed of
11 desktop computers and servers. We show that significant speedups of up to 364x compared
12 to a serial execution can be achieved using our in–house Hadoop cluster by distributing
13 the computation and automatically pruning the search space while still identifying the
14 best–performing parameter combinations. To the best of our knowledge, this is the first
15 article presenting practical results in detail for complex data analysis tasks on such a
16 heterogeneous infrastructure together with a linked simulation framework that allows for

17 computing resource planning. The results are directly applicable in many scenarios and
18 allow implementing an efficient and effective strategy for medical (image) data analysis and
19 related learning approaches.

20 **Keywords:** hyperparameter optimization; grid search; random search; support vector
21 machines; random forests; distributed computing; image analysis

22 1. Introduction

23 Exhaustive grid parameter search is a widely used hyperparameter optimization strategy in the
24 context of machine learning [1]. Typically, it is used to search through a manually defined subset
25 of hyperparameters of a learning algorithm. It is a simple tool for optimizing the performance of
26 machine learning algorithms and can explore all regions of the defined search space if no local extrema
27 exist and the surfaces of the parameter combinations are relatively smooth. However, it involves high
28 computational costs increasing exponentially with the number of hyperparameters as one predictive
29 model needs to be constructed for each combination of parameters (and possibly for each fold of a
30 Cross-Validation (CV)). It can therefore be extremely time-consuming (taking multiple days, weeks
31 or even months of computation depending on the infrastructure available) even for learning algorithms
32 with a small number of hyperparameters, which is often the case. Random search is another approach
33 that randomly samples parameters in a defined search space. It can also be very time-consuming
34 when working with a large number of hyperparameters and a large number of sample points in the
35 search space. Random search can be more suited if highly local optimal parameter combinations exist
36 that might be missed with grid search. It is a less reproducible approach though. Fortunately, grid,
37 random and similar parameter search paradigms are typically “embarrassingly parallel”¹ problems, as
38 the computation required for building the predictive model for an individual parameter setting does not
39 depend on the others [2].

40 Distributed computing frameworks can help saving time by running independent tasks simultaneously
41 on multiple computers [3] including local hardware resources, as well as Cloud computing resources.
42 These frameworks can use Central Processing Units (CPUs), Graphical Processing Units (GPUs) (which
43 have received much attention recently, especially in the field of deep learning) or a combination of
44 both. Various paradigms for distributed computing exist: Message Passing Interface (MPI)² and
45 related projects such as Open Multi-Processing (OpenMP) are geared towards shared memory and
46 efficient multi-threading. They are well-suited for large computational problems requiring frequent
47 communication between threads (either on a single computer or over a network) and are classically
48 targeted at languages such as C, C++ or Fortran. They offer fast performance but can increase
49 the complexity of software development and require high-performance networking in order to avoid
50 bottlenecks when working with large amounts of data. Other paradigms for large-scale data processing,

¹ https://en.wikipedia.org/wiki/Embarrassingly_parallel, as of 18 February 2016

² https://en.wikipedia.org/wiki/Message_Passing_Interface, as of 18 February 2016

51 including MapReduce implementations such as Apache Hadoop³, are more aimed towards data locality,
52 fault tolerance, commodity hardware and simple programming (with a stronger link to languages such
53 as Java or Python). They are more suited for the parallelization of general computation or data
54 processing tasks, with specific tools available for different kinds of processing (for example Apache
55 Spark⁴ for in-memory processing or Apache Storm⁵ for realtime stream-based computation). All of
56 these frameworks are commonly used in medical imaging and machine learning research [3,4].

57 It is also noteworthy to mention that although hyperparameter search should be as exhaustive as
58 possible, there often exist large areas of the search domain that produce suboptimal results, therefore
59 offering opportunities to intelligently reduce the search space and computation time. In a distributed
60 setting, this can complicate the process as the pruning operation requires sharing information between
61 tasks. To this end, a distributed synchronization mechanism can be designed to allow identifying
62 parameter combinations yielding suboptimal results and subsequently cancel their execution in order
63 to further decrease the total computational time. Moreover, parameter search can be a lengthy process,
64 even when executed within a distributed environment. Therefore, the availability of a parallel execution
65 simulation tool can help estimate the total runtime for varying conditions, such as the number of
66 available computation tasks. Such a simulation tool can also be useful for price estimation when
67 using “Pay-as-you-go” computing resources in the Cloud (most Cloud providers offer specific Hadoop
68 instance types and simple cluster setup tools). This allows making a trade-off between the expected
69 optimization of parameters vs. the related costs.

70 In this article, we present a novel practical framework for the simulation, optimization and execution
71 of parallel parameter search for machine learning algorithms in the context of medical image analysis.
72 It combines all the aspects discussed above: (i) parallel execution of parameter search, (ii) intelligent
73 identification and cancellation of suboptimal parameter combinations within the distributed environment
74 and (iii) simulation of the total parallel runtime according to the number of computing nodes available
75 when executed in a distributed architecture. The objective is to allow easily running very fine-grained
76 grid or random parameter search experiments in a reasonable amount of time, while maximizing the
77 likelihood of finding one of the best-performing parameter combinations. We evaluated our framework
78 with two use-cases in the article: lung tissue identification in Computed Tomography (CT) images
79 using (I) Support Vector Machines (SVMs) based on a Radial Basis Function (RBF) kernel and (II)
80 Random Forests (RFs). Results for both grid and random search strategies are provided. The main
81 contributions of the article concern the practical design, implementation and testing of a distributed
82 parameter optimization framework, leveraging software such as Hadoop and ZooKeeper in order to
83 enable efficient distributed execution and synchronization, intelligently monitoring the global evolution
84 of the grid search and canceling poorly performing tasks based on several user-defined criteria, on
85 real data and with a real problem in a scenario potentially similar to many research groups in data
86 science. This has not been done so far, to the best of our knowledge. A second contribution is the
87 developed simulation tool that allows estimating costs and benefits for a large number of scenarios prior

³ <http://hadoop.apache.org/>, as of 18 February 2016

⁴ <http://spark.apache.org/>, as of 18 February 2016

⁵ <http://storm.apache.org/>, as of 18 February 2016

88 to choosing the solution that is optimal for specific constraints. Compared to other publications with
89 a more theoretical focus on hyperparameter optimization algorithms or system design principles, such
90 as [2,5–9], this paper describes a distributed framework which is already implemented and working
91 and has been tested on medical imaging data as an example application field. Only a small number of
92 parameters were optimized in this case but the same framework also applies to larger parameter spaces.

93 The rest of the article is structured as follows : Section 2 discusses existing projects, tools and articles
94 related to the task of hyperparameter optimization. Section 3 presents the datasets, existing tools and
95 algorithms that were used. The implementation of the developed framework and the experimental results
96 obtained are detailed in Section 4. The findings and limitations are discussed in Section 5. Finally,
97 conclusions are drawn and future work is outlined in Section 6.

98 2. Related Work

99 Extensive research has already been conducted in the field of optimizing and improving on the
100 classical grid parameter search model and achieving more efficient hyperparameter optimization in the
101 context of machine learning applications. In 2002, Chapelle *et al.* proposed a method for tuning kernel
102 parameters of SVMs using a gradient descent algorithm [10]. A method for evolutionary tuning of
103 hyperparameters in SVMs using Gaussian kernels was proposed in [7]. Bergstra *et al.* [2] showed that
104 using random search instead of a pure grid search (in the same setting) can yield equivalent or better
105 results in a fraction of the computation time. Snoek *et al.* proposed methods for performing Bayesian
106 optimization of various machine learning algorithms, which supports parallel execution on multiple
107 cores and can reach or surpass human expert-level optimization in various use-cases [9]. Bergstra *et al.*
108 *al.* also proposed novel techniques for hyperparameter optimization using a Gaussian process approach
109 in order to train neural networks and Deep Belief Networks (DBNs). They proposed the Tree-structured
110 Parzen Estimator (TPE) approach and discuss the parallelization of their techniques using GPUs [11].
111 These papers discuss more the theoretical aspects of optimization, presenting algorithms but not concrete
112 implementations on a distributed computing architecture.

113 An extension to the concept of Sequential Model-Based Optimization (SMBO) was proposed
114 in [6], allowing for general algorithm configuration in a cluster of computers. The paper's focus is
115 oriented towards the commercial CPLEX solution and not an open-source solution such as Hadoop.
116 Auto-WEKA, described in [8], goes beyond simply optimizing the hyperparameters of a given machine
117 learning method, allowing for an automatic selection of an efficient algorithm among a wide range
118 of classification approaches, including those implemented in the Waikato Environment for Knowledge
119 Analysis (WEKA) machine learning software, but no distributed architecture is discussed in the article.
120 Another noteworthy publication is the work by Luo [5], who presents the vision and design concepts
121 (but no description of the implementation) of a system aiming to enable very large-scale machine
122 learning on clinical data, using tools such as Apache Spark and its MLlib machine learning library.
123 The design includes clinical parameter extraction, feature construction and automatic model selection
124 and tuning, with the goal of allowing healthcare researchers with limited computing expertise to easily
125 build predictive models.

126 Several tools and frameworks have also been released, such as the SURrogate MOdeling (SUMO)
127 Toolbox [12] that enables model selection and hyperparameter optimization. It supports grid or cluster
128 computing but it is geared towards more traditional grid infrastructures such as the Sun/Oracle Grid
129 Engine, rather than more modern solutions such as Apache Hadoop, Apache Spark, etc. Another
130 example is Hyperopt [13], a Python library for model selection and hyperparameter optimization that
131 supports distributed execution in a cluster using MongoDB⁶ for inter-process communication, currently
132 for random search and TPE algorithms⁷. It does not take advantage of the robust task scheduling
133 and distributed storage features provided by frameworks like Apache Hadoop. In the field of scalable
134 machine learning, Apache Mahout⁸ allows running several classification algorithms (such as Random
135 Forests or Hidden Markov Models) as well as clustering algorithms (k-Means Clustering, Spectral
136 Clustering, etc.) directly on a Hadoop cluster [4], but it does not address hyperparameter optimization
137 directly and also does not currently provide implementations for certain important classification
138 algorithms such as SVMs. The MLlib machine learning library⁹ provides similar features, using the
139 Apache Spark processing engine instead of Hadoop. Sparks *et al.* describe the TuPAQ system in [14],
140 an extension of the MLbase¹⁰ platform, which is based on Apache Spark's MLlib library. TuPAQ allows
141 automatically finding and training predictive models on an Apache Spark cluster. It does not mention a
142 simulation tool that could help estimating the costs of running experiments of varying complexity in a
143 Cloud environment.

144 Regarding the early termination of unpromising results (pruning the search space of a parameter
145 search) in a distributed setting, [15] describes a distributed learning method using the multi-armed bandit
146 approach with multiple players. SMBO can also incorporate criteria based on multi-armed bandits [11].
147 This is also related to the early termination approaches proposed in this paper that are based on the first
148 experiments and cutoff parameters based on our experiences.

149 However, articles describing a distributed parameter search setup in detail, including the framework
150 used and an evaluation with real-world clinical data, are scarce. A previous experiment on a much
151 smaller scale was conducted in [3], where various medical imaging use-cases were analyzed and
152 accelerated using Hadoop. A more naive termination clause was used in a similar SVM optimization
153 problem, where suboptimal tasks were canceled based on a single decision taken after processing a fixed
154 number of patients for each parameter combination, based solely on a reference time set by the fastest
155 task reaching the given milestone. The approach taken in this paper is more advanced and flexible, as
156 it cancels tasks during the whole duration of the job, based on an evolving reference value set by all
157 running tasks.

158 In this article we describe a very practical approach in detail, based on the Hadoop framework that
159 is easy to set up and manage in a small computing environment, but also easily scalable for larger

⁶ <http://mongodb.org/>, as of 18 February 2016

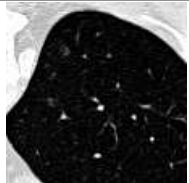
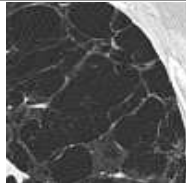

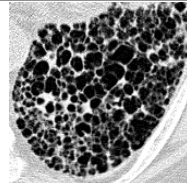
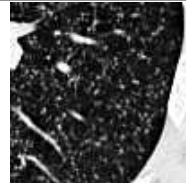
⁷ <http://jaberg.github.io/hyperopt/>, as of 18 February 2016

⁸ <http://mahout.apache.org>, as of 18 February 2016

⁹ <http://spark.apache.org/mllib/>, as of 18 February 2016

¹⁰ <http://mlbase.org>

Table 1. Visual aspect and distribution of the 32×32 blocks per class of lung tissue pattern. A patient may have several types of lung disorders.

| visual aspect |  |  |  |  |  |
|-----------------------|---|---|---|---|---|
| tissue type | healthy | emphysema | ground glass | fibrosis | micronodules |
| hand-drawn ROIs | 150 | 101 | 427 | 473 | 297 |
| 32×32 blocks | 5167 | 1127 | 2313 | 3113 | 6133 |
| patients | 7 | 6 | 32 | 37 | 16 |

160 experiments and supported by many Cloud infrastructure providers if the locally available resources
161 become insufficient.

162 3. Material and Methods

163 This section describes the datasets, tools and experimental setup used for developing and testing
164 the parallel parameter search framework. It also details the testing use-cases used to evaluate the
165 framework and the adaptive criteria for canceling tasks corresponding to parameter combinations leading
166 to suboptimal classification performance.

167 3.1. Datasets

168 The medical image classification task used for this article consists of the identification of five lung
169 texture patterns associated with interstitial lung diseases in high-resolution CT images [16]. The image
170 instances consist of 2D 32×32 blocks represented in terms of the energies of sixth-order aligned Riesz
171 wavelet coefficients [17,18], yielding a feature space with 59 dimensions when concatenated with 23
172 intensity-based features. The distribution and visual aspect of the lung tissue types (including the number
173 of hand-drawn Regions of Interest (ROIs), blocks and patients) are detailed in Table 1. Going towards
174 full 3D data analysis also increases runtime for this use case even more but the current data with larger
175 inter-slice distance does not allow for this.

176 3.2. Existing Tools

177 The developed framework relied on Apache Hadoop¹¹ and can be used with any kind of parameter
178 search problem. Hadoop is a distributed storage and computation tool that supports the MapReduce
179 programming model made popular by Google [19] (among others, such as Apache Spark or Apache
180 Storm). Use of Hadoop is frequent in medium-sized research groups in data science, as it is quick and
181 easy to set up and use, also on heterogeneous infrastructures.

¹¹ <http://hadoop.apache.org/>, as of 24 February 2016

182 The MapReduce model is used in the context of our experiments, as it is simple and fits our needs
183 well. It separates large tasks into 2 phases, called “Map” and “Reduce”. In a typical setting, the “Map”
184 phase splits a set of input data into multiple parts, which are further processed in parallel and produce
185 intermediate outputs. The “Reduce” phase aggregates the intermediate outputs to produce the final job
186 result. In the context of this article, we only implemented the “Map” phase, as no aggregation was
187 required on the output of this first phase.

188 Hadoop consists of two main components. The first is a distributed data storage system called Hadoop
189 Distributed File System (HDFS) that manages the storage of extremely large files in a distributed,
190 reliable and fault-tolerant manner. It was used for data input and output when running computations.
191 A detailed description of HDFS can be found in [20]. The second component is the distributed data
192 processing system that was called Hadoop MapReduce in early versions of the software and Yet Another
193 Resource Negotiator (YARN) since version 2.0 of Hadoop. The reason behind the name change is
194 that the programming algorithm was decoupled from the execution framework in the second generation
195 of Hadoop, allowing for more flexible use of different distributed programming paradigms, i.e., it
196 is not restricted to the batch-oriented MapReduce framework [21] anymore. This can also provide
197 opportunities for making the developed framework evolve towards new paradigms and use-cases.

198 The synchronization of distributed parallel tasks was performed with Apache ZooKeeper¹². The focus
199 of this tool is to provide highly reliable distributed coordination [22]. The architecture of ZooKeeper
200 supports redundancy and can therefore provide high availability. The data are stored in the computation
201 nodes and are saved under hierarchical name spaces, similar to a file system or other tree structures.

202 The simulation tool used for estimating the runtime of a Hadoop job under given conditions (as well
203 as tweaking parameters of the experiments) was programmed in Java and is detailed in Section 4.2. It
204 uses the output of one full Hadoop job as a baseline for running simulations. The WEKA Data Mining
205 Software [23] was used for the implementation of the SVM and RF classifiers.

206 3.3. Hardware and Hadoop Cluster

207 The in-house Hadoop cluster consisted of:

- 208 • 21 nodes including a majority of 8-core CPU desktop stations with 16 Gigabytes (GBs) of
209 Random-Access Memory (RAM), as well as 4 more powerful machines (24 cores and 64GB
210 of RAM, 24 cores and 96GB of RAM, 40 cores and 128GB of RAM, 64 cores and 128GB of
211 RAM).
- 212 • Gigabit Ethernet network connections between all nodes.
- 213 • A total of 152 simultaneous Map tasks (number of cores attributed to Hadoop in the cluster) and
214 26 simultaneous Reduce tasks. The total is given by the number of tasks that were assigned to the
215 Hadoop cluster on each node, both for the Map and Reduce phases.

¹² <http://zookeeper.apache.org/>, as of 24 February 2016

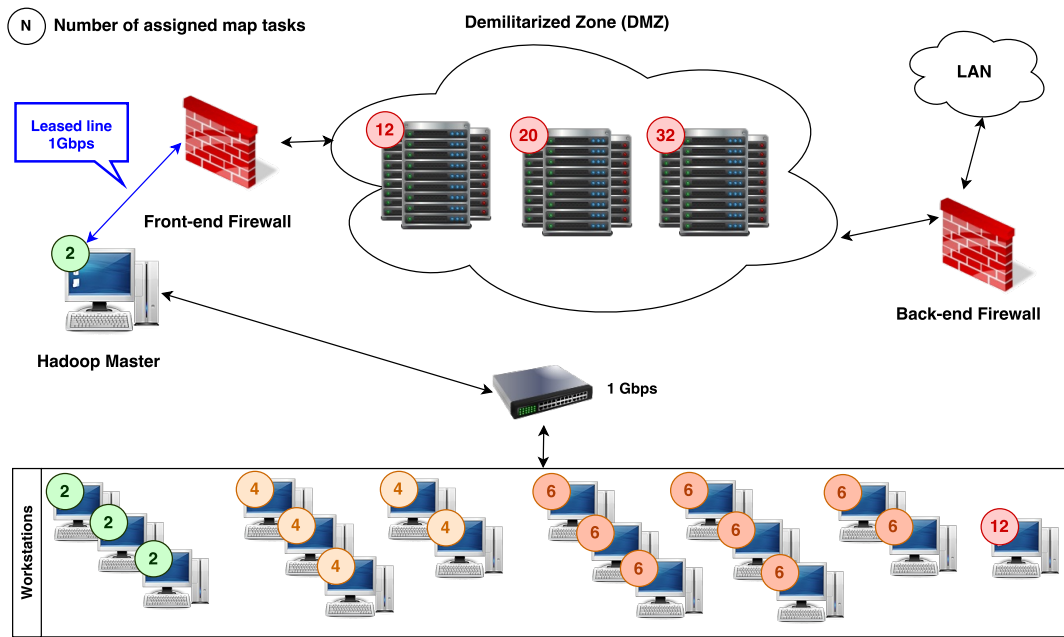


Figure 1. Schema of the in-house Hadoop cluster, showing all the nodes and the number of assigned Map tasks.

216 Figure 1 shows a schema of the cluster of machines, listing all the nodes and the network configuration, as
 217 well as the number of Map tasks assigned to each computer, as nodes are configured according to their
 218 computing power. All desktop machines are commonly used by researchers during the day, therefore
 219 only a subset (usually about 50%) of CPU cores and main memory are attributed to the Hadoop cluster.
 220 Previous research showed that the daily normal usage of machines has little impact on the duration of
 221 Hadoop jobs in our environment [3].

222 3.4. Classification Algorithms

223 Two classification algorithms were used and optimized for the categorization of the lung tissue types:
 224 SVMs and RFs. An extension to other tasks is easily possible but these two are characteristic for many
 225 other techniques and both are frequently used in machine learning and medical imaging.

226 SVMs have shown to be effective to categorize texture in wavelet feature spaces [24] and in particular
 227 for lung tissue [25]. Kernel SVMs implicitly map feature vectors \mathbf{v}_i to a higher-dimensional space by
 228 using a kernel function $K(\mathbf{v}_i, \mathbf{v}_j)$. We used the RBF kernel given by the multidimensional Gaussian
 229 function

$$K(\mathbf{v}_i, \mathbf{v}_j) = e^{\frac{-\|\mathbf{v}_i - \mathbf{v}_j\|^2}{2\gamma}}. \quad (1)$$

230 SVMs build separating hyperplanes in the higher-dimensional space considering a two-class
 231 problem. Two parallel hyperplanes are constructed symmetrically on each side of the hyperplane that
 232 separates the two classes. The goal of SVMs is to maximize the distance between the two external
 233 hyperplanes, called the margin [26]. This yields the decision function $f(\mathbf{v}_i)$, which minimizes the
 234 functional

$$\|f\|_K + C \sum_{i=1}^N \max(0, 1 - y_i f(\mathbf{v}_i))^2, \quad (2)$$

with $\|f\|_K$ the norm of the reproducing kernel Hilbert space defined by the kernel function K , N the total number of feature vectors, and y_i the class labels (i.e., $y_i \in \{-1, 1\}$). The parameter C determines the cost attributed to errors and requires optimization to tune the bias–variance trade–off. For multiclass classification, several one–versus–all classifiers are built and the model with the highest decision function determines the predicted class. Two parameters are being optimized for SVMs: the cost C and the parameter of the Gaussian kernel γ .

RFs consist of building ensembles of Decision Trees (DTs) [27]. Each DT is built on a subset of features and a subset of training vectors (i.e., bagging). The DTs divide the feature space successively by choosing primarily features with the highest information gain [28]. The final class prediction of RFs is obtained as the mean prediction of all individual trees. Three parameters are being optimized for RFs: the number of generated random trees T , and for each DT: the number of randomly selected features F and the maximum tree depth D .

3.5. Task Cancellation Criteria

The following is a description of the method used for deciding which hyperparameter combinations to keep during the execution of the experiments on the Hadoop cluster. The classification accuracy acc_k associated with one set of hyperparameters is monitored throughout the execution of the k folds of the CV. In order to determine if a hyperparameter combination is performing well, the first considered criterion is whether the value of acc appears to be stable for the given combination over the k folds of the CV. The mean accuracy μ_{acc} is updated each time a new value for this hyperparameter combination is available (i.e., each time that a new fold of the CV has completed) and added to a list of values. At the same time, the variance σ_{acc} is calculated for the set of recorded mean accuracies over a “sliding window” of size W_k . Finally, the gradient of the variance is determined over these W_k values as $\frac{\partial \sigma_{acc}}{\partial k}$, $k \in [1, \dots, W_k]$. $\frac{\partial \sigma_{acc}}{\partial k}$ was computed using least squares regression. If the gradient is $\frac{\partial \sigma_{acc}}{\partial k} \leq 0$, the estimated classification accuracy was considered to be stable, otherwise the evolution of the mean accuracy is deemed to be unstable and no decision is taken yet about the cancellation of this combination. When the accuracy is found to be stable, the second step consists of comparing one or more criteria of the current combination of parameters against the global evolution of the classification accuracy given by all other parameter combinations. Two criteria are considered:

- Is the current mean accuracy of the combination μ_{acc} lower than the global mean accuracy (minus a margin of Δ_{acc})?
- Is the current mean runtime of 1 task for the combination longer than the global mean runtime for 1 task (multiplied by a factor of Δ_t)?

The first criterion is monitoring the accuracy of the current hyperparameter combination. Given that σ_{acc} is considered to be stable, the chance that the accuracy associated with this combination of parameters improves significantly later is relatively small. Therefore, the combination is canceled if its current

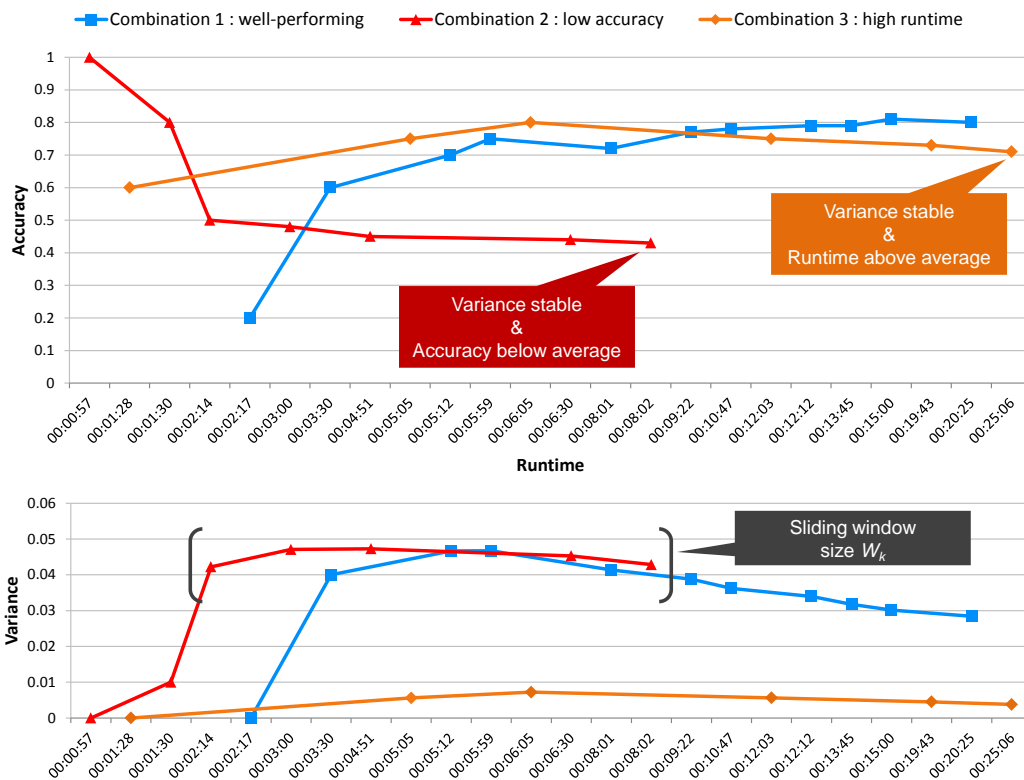


Figure 2. Illustration of the task cancellation process for SVMs. The top graph shows the evolution of the mean accuracy μ_{acc} and the bottom graph plots the evolution of the variance thereof for 3 parameter combinations (well-performing, low accuracy and high runtime). The cancellation checks are performed for each combination only when at least W_k variance values are available and the evolution of the variance is considered to be stable (see Section 3.5).

270 accuracy is lower than the current global accuracy. The second criterion works in a similar fashion but
 271 is based on the runtime of the tasks. Indeed, for certain classifiers such as SVMs the longer the time to
 272 achieve convergence, the higher the likelihood of a bad performance [3]. For this reason, abnormally
 273 time-consuming parameter combinations are also canceled, because they generally yield suboptimal
 274 results and more importantly have a significant impact (as much as one order of magnitude higher
 275 than average runtimes) on the overall runtime of the experiment when not canceled. Δ_{acc} and Δ_t can
 276 be tuned to balance between overall computational time and classification performance. Additionally,
 277 each criterion can be individually enabled or disabled, as not all classifiers follow the same behavior.
 278 Algorithm 1 outlines the process described above, with current mean values for the accuracy and runtime
 279 being obtained first (both global and for the given parameter combination $pComb$), followed by the set
 280 of variances of size W_k . Subsequently, the stability test described in this section is performed, as well
 281 as the performance checks (accuracy and runtime) in case of a stable evolution. If the combination is
 282 performing poorly, its status is set to 'cancelled'. Figure 2 shows an illustration of how the cancellation
 283 process works with 3 parameter combinations: a well-performing combination, a combination with
 284 suboptimal accuracy and a combination with above-average runtime.

Algorithm 1 Parameter combination cancellation

```

1: function CANCELPARAMETER( $pComb$ )
2:    $\mu_{accGlobal} \leftarrow \text{currGlobalMeanAcc}()$ 
3:    $\mu_{acc} \leftarrow \text{currMeanAcc}(pComb)$ 
4:    $\mu_{timeGlobal} \leftarrow \text{currGlobalMeanRuntime}()$ 
5:    $\mu_{time} \leftarrow \text{currMeanRuntime}(pComb)$ 
6:    $setOf\sigma_{acc} \leftarrow \text{currVarsOfMeanAccs}(pComb)$ 
7:   if varIsStable( $setOf\sigma_{acc}$ ) then
8:     if  $\mu_{acc} < \mu_{accGlobal} - \Delta_{acc}$  or  $\mu_{time} > \mu_{timeGlobal} - \Delta_t$  then
9:        $pComb.status \leftarrow \text{'cancelled'}$ 
10:    end if
11:  end if
12: end function

```

285 **4. Results**

286 This section describes the implementation of the framework and experimental results obtained.

287 *4.1. Implementation of the Hadoop-based Execution Framework*288 *4.1.1. Standard Run*

289 The following list outlines all the chronological steps for running a distributed parameter search using
 290 the framework, but without optimization (i.e., no task cancellation). This is referred to as the **standard**
 291 **run**.

- 292 1. An input file containing a hash table with all the possible combinations of parameter values and
 293 patient identifiers (the latter was used for performing a Leave–One–Patient–Out (LOPO) CV) is
 294 created (one combination per line). This hash table was based on parameter ranges specified by
 295 the user. In the case of a random search, the user simply specifies the lower and upper bounds of
 296 each parameter, the values are then generated randomly within this space. The order of the lines
 297 was randomized in order to avoid executing a large number of similarly complex tasks at the same
 298 time. The file is then uploaded to the HDFS where it serves as the input file of the Hadoop job.
- 299 2. The Hadoop job starts, splitting the workload into N/M Map tasks, where N is the total number
 300 of lines in the file and M is a variable defining how many lines a single task should process. M
 301 can be tweaked in order to avoid having Map tasks that are extremely short (less than 10 seconds).
 302 Map tasks that are too short can impact the runtime of a Hadoop job in a non–negligible fashion
 303 due to overhead caused by starting and managing Hadoop tasks.
- 304 3. Each task executes a setup function (only once per Map task) that contains the following steps:
 305 (a) Load the dataset and prepare it for use (in this case, set the instance class attribute).
 306 (b) Normalize the dataset: the feature values were scaled to $[0, 1]$.

- 307 4. Each task executes the Map function (M times per task) that consists of one fold of the LOPO CV:
- 308 (a) Split the data into a **training** set containing all the instances of the dataset except for those
- 309 of the current patient and a **testing** set containing all the instances of the current patient.
- 310 (b) Build the classifier using the current combination of parameters (for example C and γ in the
- 311 SVM use–case) and the training set.
- 312 (c) Classify each instance of the test set using the previously built classifier model.
- 313 (d) Get the number of total and correctly classified instances and write them as the output of the
- 314 function.

315 The above process is shown in Algorithm 2.

Algorithm 2 Execution Framework - Standard Run

```

1: generateInput()
2: startJob()
3: for all  $task \in N/M_{tasks}$  do
4:   loadDataset()       $\rightarrow$  Setup ( $1 \times$  per task)
5:   prepareDataset()
6:   normalizeDataset()
7:   for all  $pComb \in M_{pCombinations}$  do       $\rightarrow$  Map ( $M \times$  per task)
8:      $tStart \leftarrow \text{NOW}$ 
9:      $classifierArgs \leftarrow pComb.classifierArgs$ 
10:     $patientID \leftarrow pComb.patientID$ 
11:    configureClassifier( $classifierArgs$ )
12:    splitDataSet( $patientID$ )       $\rightarrow$  LOPO
13:    trainClassifier( $trainingSet$ )
14:     $result \leftarrow \text{classifyTestSet}(testSet)$ 
15:     $tEnd \leftarrow \text{NOW}$ 
16:     $acc \leftarrow result.correct/result.total$ 
17:     $runtime \leftarrow tEnd - tStart$ 
18:    writeOutput( $acc, result.correct, result.total, patientID, runtime$ )
19:   end for
20: end for

```

316 4.1.2. Optimized Run

317 When activating the mode that cancels suboptimal tasks (referred to as **optimized run**, see Section

318 3.5), the process was slightly modified :

- 319 1. Before the job starts, various “znodes” (i.e., znodes are files persisted in memory on the ZooKeeper
- 320 server) are initialized for storing parameter combination accuracy values, a list of canceled
- 321 parameter combinations, etc.

- 322 2. During the setup (point 3 of the previous list), a connection to the ZooKeeper object is established
323 and the variables for canceling tasks are attributed.
- 324 3. At the start of the Map function (point 4 of the previous list), a check is performed to identify if the
325 parameter combination was already canceled. If this is the case, the function returns immediately,
326 otherwise the classification is performed as usual.
- 327 4. Once the classification is finished, several values in the ZooKeeper server are updated :
- 328 (a) The number of total and correctly classified instances, as well as the runtime of the tasks for
329 the given parameter combination are incremented.
- 330 (b) The current accuracy of the given parameter combination was added to the
331 `DescriptiveStatistics` object (part of the Apache Commons Math library¹³), which
332 allows easy calculations of statistical values (e.g., μ_{acc} , σ_{acc}) on an evolving set of data.
- 333 (c) The current variance σ_{acc} (computed from all existing accuracies for the given parameter
334 combination) is added to a circular buffer of size W_k . This buffer is further used to calculate
335 the gradient of the variance evolution over the last W_k values.
- 336 (d) The number of total and correctly classified instances, as well as the runtime of the tasks for
337 the global job are incremented.
- 338 5. At the end of the Map function, a check is performed whether the current parameter combination
339 needs to be canceled or not. This check takes into account the following variables:
- 340 (a) Variance over the last W_k values is stable, i.e. $\frac{\partial \sigma_{acc}}{\partial k} \leq 0$. If the gradient is positive, it is
341 assumed that the values are still changing significantly and the parameter combination is not
342 canceled.
- 343 (b) Mean accuracy of the given parameter combination. If μ_{acc} is smaller than the mean global
344 accuracy of all parameter combinations minus a Δ_{acc} (set to 0.05 in our experiments),
345 the parameter combination is canceled (or blacklisted), i.e. the classification step in all
346 subsequent Map tasks of the corresponding parameter combination will not be executed.
- 347 (c) Mean runtime of the given parameter combination. If the latter is longer than the mean global
348 runtime of all parameter combinations multiplied by a Δ_t (set to 2.0 in our experiments), the
349 parameter combination is canceled, i.e. the classification step in subsequent Map tasks of the
350 corresponding parameter combination is not executed.

351 The above process is shown in Algorithm 3, where differences with the standard run (Algorithm 2) are
352 highlighted.

¹³ <http://commons.apache.org/proper/commons-math/>, as of 24 February 2016

Algorithm 3 Execution Framework - Optimized Run (differences with Algorithm 2 are highlighted)

```

1: generateInput()
2: initSyncFields()
3: startJob()
4: for all  $task \in N/M_{tasks}$  do
5:   loadDataset()       $\rightarrow$  Setup ( $1 \times$  per task)
6:   prepareDataset()
7:   normalizeDataset()
8:   initStatsObjects()
9:   for all  $pComb \in M_{pCombinations}$  do       $\rightarrow$  Map ( $M \times$  per task)
10:  if  $pComb.status = 'cancelled'$  then
11:    continue       $\rightarrow$  Skip cancelled iterations
12:  end if
13:   $tStart \leftarrow \text{NOW}$ 
14:   $classifierArgs \leftarrow pComb.classifierArgs$ 
15:   $patientID \leftarrow pComb.patientID$ 
16:  configureClassifier( $classifierArgs$ )
17:  splitDataSet( $patientID$ )       $\rightarrow$  LOPO
18:  trainClassifier( $trainingSet$ )
19:   $result \leftarrow \text{classifyTestSet}(testSet)$ 
20:   $tEnd \leftarrow \text{NOW}$ 
21:   $acc \leftarrow result.correct/result.total$ 
22:   $runtime \leftarrow tEnd - tStart$ 
23:  writeOutput( $acc, result.correct, result.total, patientID, runtime$ )
24:  updateSyncFields( $pComb, result.correct, result.total, runtime$ )
25:  cancelParameter( $pComb$ )
26:  end for
27: end for

```

353 *4.2. Implementation of the Simulation Tool*

354 Time is often a limiting factor when running experiments, and it can have a strong influence on
355 the achieved results. Having a tool that can run simulated grid search experiments (modeled after
356 the real-world Hadoop-based framework) in a single machine in order to approximate runtime and
357 give indications about the expected performance can help in designing experiments, choosing sensible
358 margins for parameter cancellation (see Section 3.5), estimating the required scale of computation
359 cluster, as well as calculating the cost of running the experiment in a cloud-based "Pay-as-you-go"
360 platform. This section details the implementation of this tool: the behavior of the real-world Hadoop
361 implementation was closely reproduced, with the following characteristics and differences:

- 362 • The results of a Hadoop experiment (containing the runtime of each task) are loaded into the
363 simulator Java class: they will serve as a baseline for simulating Hadoop jobs with different
364 amounts of available computation tasks and different values for the termination criteria margins,
365 for example.
- 366 • A "time step" counter is initialized and incremented in milliseconds, simulating the passage of
367 time.
- 368 • A queue of running tasks (of size T , representing the number of Map tasks in the simulated cluster)
369 is populated.
- 370 • After each millisecond, the starting / ending tasks are managed and the cancellation checks are
371 performed like on the Hadoop cluster. The major difference is that instead of using the ZooKeeper
372 distributed synchronization system, simple Java data structures are used (hash maps, lists, etc.) for
373 monitoring the evolution of parameter combination performance.
- 374 • Each time a task completes, another pending task is added to the queue of running tasks. This
375 behavior is the same as in Hadoop.
- 376 • At the end of the simulated Hadoop job (all tasks are processed), statistics about the simulated
377 job are given as an output: total duration of the job (if executed in a real cluster of a given size),
378 number of canceled parameters, maximum achieved accuracy, etc.

379 The goal is to have a tool that can provide an approximation of the average runtime of a task for a given
380 machine learning scenario, including the variance in processing time for different parameter values. A
381 real small-scale experiment with a coarse grid can be run to get a clear idea of these values, that can
382 then serve as a base for a simulated experiment at a much larger scale. If running a real experiment
383 before simulation is not desired or feasible, the tool can also easily use an average runtime per task (with
384 margins to represent shorter and longer tasks) directly input by the user after performing some local
385 empirical tests.

386 *4.3. Experimental Results*

387 Several experiments were performed:

- 388 • Determining the speedups that can be obtained (with and without task cancellation) compared to
389 a serial execution on a single computer.
- 390 • Verifying whether the best parameter combination is kept when canceling tasks,
- 391 • Comparing the runtime and performance between grid and random search,
- 392 • Investigating if the developed simulation tool can provide a realistic approximation of the runtime
393 of an experiment under varying conditions.

394 *4.3.1. Grid Search*

395 The first experiments were conducted with the classical grid parameter search strategy. All
396 the experiments were run using the Hadoop cluster configuration described in Section 3.3. When
397 the objective function (e.g., classification accuracy) is expected to be smooth through consecutive
398 parameters, the grid search is expected to lead to reproducibly good results with a trade-off between
399 grid size and the probability to find the maximum performance (or be at least very close to it).

400 For both use-cases (RF and SVM), the Hadoop job was run twice : once based on the **standard run**
401 mode, where no tasks were canceled during the execution of the job and once based on the **optimized**
402 **run** mode, where tasks corresponding to suboptimal parameter combinations were canceled. The results
403 are presented in Table 2. An estimation of the time required to run the computation serially on a single
404 computer is provided in the first column. The estimation is based on the runtime recorded for each Map
405 task, purely for the classification part, therefore excluding the overhead produced by Hadoop for starting
406 and managing tasks.

407 The grid parameter search domain is defined as follows:

- 408 • For the SVM use-case, two parameters are being optimized.
 - 409 1. The cost C , varying from 0 to 100 in increments of 10 (and $C = 0$ is replaced with $C = 1$).
 - 410 2. The kernel parameter G , varying from -2.0 to 2.0 in increments of 0.1 (actual kernel value
411 is computed by $\gamma = 10^G$)
- 412 • For the RF use-case, three parameters are being optimized.
 - 413 1. The number of trees in the random forest T , varying from 0 to 1000 in increments of 10 (and
414 $T = 0$ is replaced with $T = 10$)
 - 415 2. The maximum tree depth D , varying from 0 to 4 in increments of 1 (where $D = 0$ signifies
416 that the depth is not restricted)
 - 417 3. The number of randomly selected features F for testing at each node, varying from 1 to
418 $2 * \sqrt{N}$ in increments of 1 (where N is the total number of features, in this case 58)

Table 2. Experimental results showing the comparison between an estimation of running the grid parameter search on a single computer and running it on the in-house Hadoop cluster in the **standard run** and **optimized run** configurations, for both use-cases (RF and SVM). The indication in brackets [...] for the “Best accuracy” value in the **optimized run** column shows whether the best or second best achieved accuracy of the **standard run** was kept running.

| RF Optimization | | | |
|--|---------------------------------|----------------------------------|-----------------------------------|
| | Single computer (estimation) | Hadoop cluster (standard run) | Hadoop cluster (optimized run) |
| Job Execution time | 302d 00h 15m 24s | 51h 27m 03s | 19h 52m 05s |
| Total combinations | 651,460 | 651,460 | 651,460 |
| Lines per task | N/A | 20 | 20 |
| Total Map tasks | N/A | 32,573 | 32,573 |
| Number of canceled parameter combinations | N/A | 0 | 5052 / 7500 |
| Best accuracy | 0.73762 | 0.73762 | 0.73756 [2 nd BEST] |
| Speedup | 1x | ~141x | ~364x |
| SVM Optimization | | | |
| | Single computer (estimation) | Hadoop cluster (standard run) | Hadoop cluster (optimized run) |
| Job Execution time | 52d 17h 27m 10s | 8h 49m 51s | 4h 40m 57s |
| Total combinations | 38,786 | 38,786 | 38,786 |
| Lines per task | N/A | 2 | 2 |
| Total Map tasks | N/A | 19,393 | 19,393 |
| Number of canceled parameter combinations | N/A | 0 | 236 / 451 |
| Best accuracy | 0.77999 | 0.77999 | 0.77999 [BEST] |
| Speedup | 1x | ~143x | ~270x |

Table 3. Comparison between running a grid parameter search and a random search (with the same number of combinations), with and without task cancellation, for optimizing the hyper-parameters of the SVM experiment.

| Grid search and random search comparison — SVM experiment | | | | |
|--|-------------------------------|---------------------------------|--------------------------------|----------------------------------|
| | Grid search (standard run) | Random search (standard run) | Grid search (optimized run) | Random search (optimized run) |
| Job Execution time | 8h 49m 51s | 8h 14m 32s | 4h 40m 57s | 4h 08m 44s |
| Number of canceled parameter combinations | 0 | 0 | 236 / 451 | 254/451 |
| Best accuracy | 0.77999 | 0.78033 | 0.77999 | 0.78045 |

419 4.3.2. Random Search

420 Two more experiments were run using the random search strategy for the SVM use-case in order
 421 to demonstrate the flexibility of the developed framework and investigate the possible improvements in
 422 terms of runtime and maximum achieved classification accuracy. Very good and efficient results were
 423 reported for Random Search in the past despite the fact that the results are not necessarily reproducible
 424 and thus non-optimal results are a risk, albeit with low probability [2]. In order to allow fair comparisons
 425 with grid search, the same number of points used were generated randomly in the search space based on
 426 a uniform distribution, using the same upper and lower bounds. The comparison of the results is shown
 427 in Table 3. The runtimes and results are in this case very close to the ones obtained with the grid search.

428 A series of random search experiments using a varying number of randomly sampled points were
 429 also conducted, in order to analyze the evolution of both the runtime of the Hadoop job as well as the
 430 maximum achieved accuracy. The results are shown in Figure 3.

431 Finally, multiple iterations (20 in total) of the same random search experiment (using 25% of the
 432 original number of points, i.e. 112 combinations) were run in order to determine the Relative Standard
 433 Deviation (RSD) of the maximum accuracy obtained as well as the job runtime. The results are shown
 434 in Figure 4.

435 4.4. Simulation results and validation

436 Once the output of a **standard run** was available, it was fed into the simulation tool to estimate the
 437 runtime of the same job under different conditions. For instance, the number of simultaneous Map tasks
 438 can be increased to approximate the runtime on a larger Hadoop cluster. Similarly, the Δ_{acc} and Δ_t task
 439 cancellation margins can be adjusted to evaluate the time-performance trade-off (i.e., smaller margins
 440 will lead to faster runtimes but increase the risk of canceling optimal parameter combinations).

441 To validate whether the simulation tool can produce realistic results, the SVM grid search use-case
 442 was executed four times in the Hadoop cluster:

- 443 • **standard run** and **optimized run** with 152 Map tasks, $\Delta_{acc} = 0.05$ and $\Delta_t = 2.0$
- 444 • **standard run** and **optimized run** with 64 Map tasks, $\Delta_{acc} = 0.05$ and $\Delta_t = 2.0$

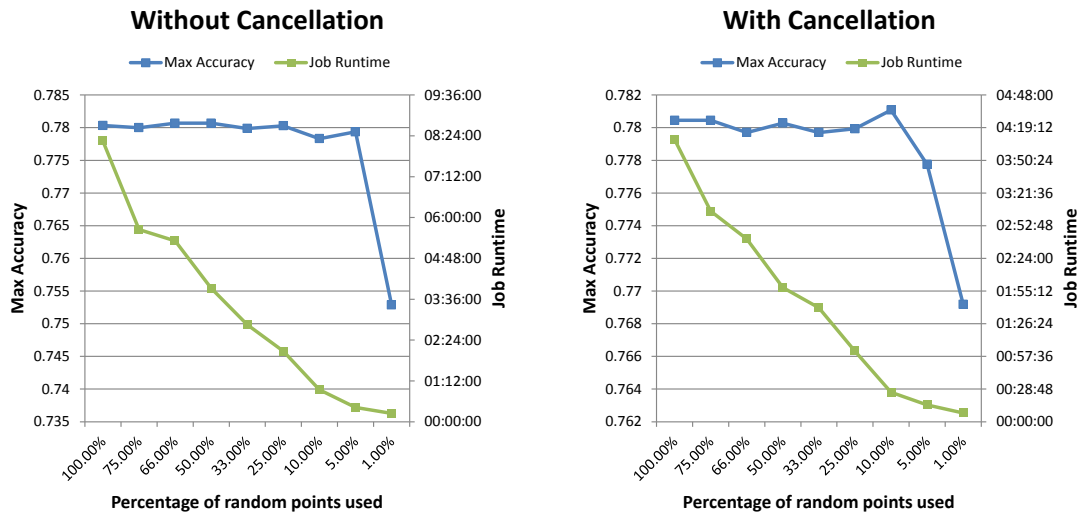


Figure 3. The graphs display the evolution of the maximum obtained accuracy and the total runtime of the SVM random search experiment (in both the **standard** and **optimized** run configurations) for a shrinking number of randomly selected points in the search space of the hyperparameters. 100% is equivalent to all 451 parameter combinations used in the comparison with the grid search method, 75% corresponds to 338 combinations, etc.

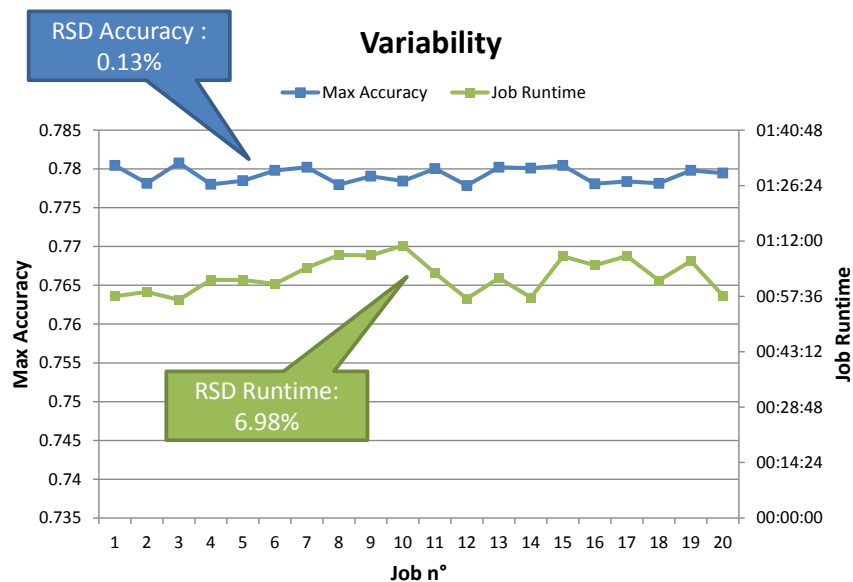
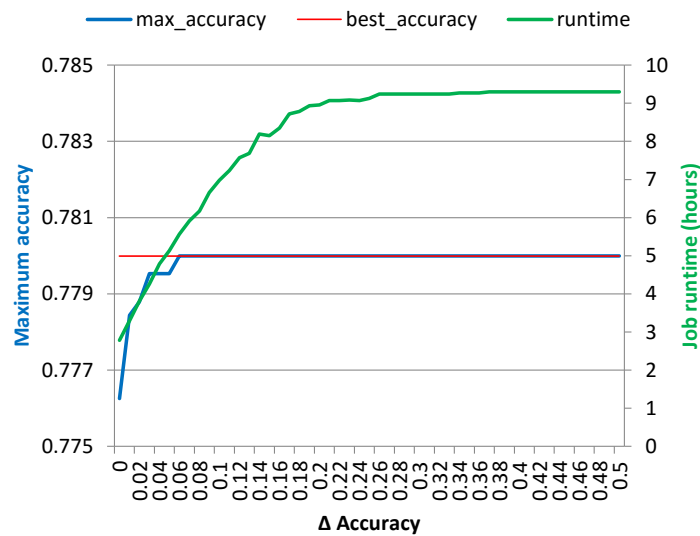


Figure 4. The graph shows the variability of the maximum obtained accuracy and the total runtime of the SVM random search experiment in the **optimized** run configuration, using 25% of the original 451 parameter combinations used in the comparison between grid and random search. The Relative Standard Deviation (RSD) of the maximum accuracy is 0.13% and the RSD of the job runtime is 6.98%.

Table 4. Validation of the simulation tool.

| SVM - standard run | | | |
|----------------------------|---------------------|-------------------|---------------------|
| | Experimental result | Simulation result | Relative difference |
| 152 Map tasks (baseline) | 9h 26m 04s | 9h 17m 56s | -1.43% |
| 64 Map tasks | 19h 42m 31s | 21h 39m 57s | +9.93% |
| SVM - optimized run | | | |
| | Experimental result | Simulation result | Relative difference |
| 152 Map tasks (baseline) | 4h 40m 57s | 5h 16m 19s | +12.58% |
| 64 Map tasks | 9h 24m 54s | 12h 17m 49s | +30.61% |

**Figure 5.** Graph displaying the evolution of the maximum obtained accuracy and the total runtime of the SVM grid search experiment for a growing margin Δ_{acc} .

445 The results of the first two executions with 152 tasks are used as the input for running four simulations
 446 with the same number of Map tasks as the runs listed above. The results are shown in Table 4. A
 447 fixed 3-second overhead (determined by empirical tests) was added to the runtime of each task in order
 448 to simulate the impact of Hadoop task setup. Moreover, an interesting opportunity provided by the
 449 simulation tool is to evaluate the effect of the cancellation margins Δ_{acc} and Δ_t on the maximum
 450 accuracy achieved. By gradually changing these values for the results of the SVM grid search
 451 experiment, Figures 5 and 6 are created. Results with cancellation are less precise in the simulations
 452 compared to results without task cancellations.

453 5. Discussion

454 Three major observations can be deduced from the experimental results: first of all, the speedup
 455 achieved by simply distributing a grid parameter search is very substantial, with the total runtime for
 456 the search accelerated by a factor of 141x (RF) and 143x (SVM), when compared to an estimation of a
 457 serial execution on a single computer (see Table 2). It also shows that the total runtime decreases almost
 458 linearly as the number of nodes (and therefore available Map tasks) in the Hadoop cluster increases.

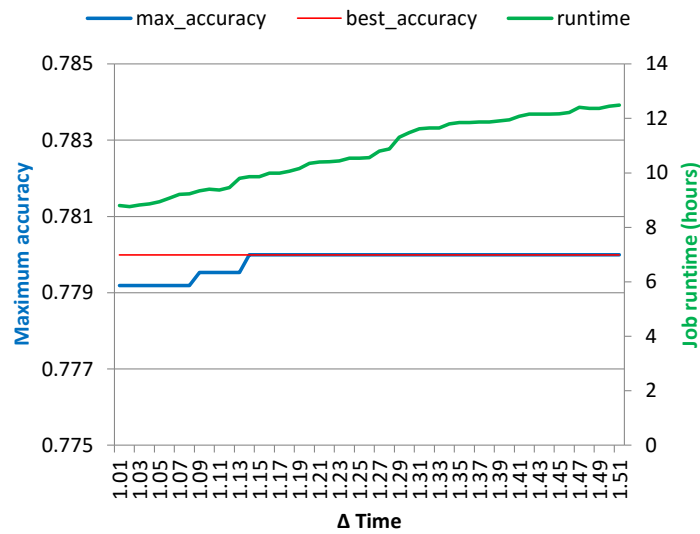


Figure 6. Graph displaying the evolution of the maximum obtained accuracy and the total runtime of the SVM grid search experiment for a growing factor Δ_t .

459 Second, adding the accuracy and runtime check and canceling suboptimal parameter combinations
 460 allows decreasing the runtime even further, by a factor close to or greater than 2x in both use-cases
 461 without any significant impact on the maximum achieved accuracy. It also shows that the framework
 462 performed well for two different types of classifiers and with a different number of hyperparameters.
 463 Third, the results show that several parameter search strategies are supported and work well with the
 464 developed framework. The random search experiments ran slightly faster than the grid search using the
 465 same number of points and gave equivalent results both with and without task cancellation (see Table 3).
 466 Moreover, reducing the number of random points yielded equivalent results in a fraction of the time
 467 needed for the grid search experiments. Repeated experiments also showed that the variability in terms
 468 of runtime and achieved performance is minimal. Random search thus provides an interesting option,
 469 also in the simulation tool.

470 The proposed simulation tool was successfully used to estimate job runtimes using a varying number
 471 of tasks, with a relative difference of ~10% between the real-world experiment and the simulation for
 472 the **standard run** using a smaller number of simultaneous tasks (64, see Table 4). For the **optimized**
 473 **run**, the errors were larger, about ~12.5% when simulating with the original amount of Map tasks,
 474 and ~30.6% when using the smaller number of tasks. Moreover, the simulation provided insights into
 475 the effect of varying the cancellation conditions on the maximum achieved classification accuracy and
 476 overall job runtime without requiring to run a battery of lengthy Hadoop jobs. The latter can be used to
 477 reduce costs when using “Pay-as-you-go” computing resources in the cloud, which might in the future
 478 become the main computation source for many research departments in any case.

479 Some limitations of this work include the LOPO CV, which could benefit from an added inner
 480 Cross-Validation (CV) performed on the training set, in order to reduce the risk of overfitting.
 481 Fortunately this is entirely possible with our framework and is well-suited for parallelizing the task
 482 even further. Another limitation concerns the simulation tool, which currently works based only on the
 483 results of an real-world experiment. Although it is still interesting to use it on a small-scale experiment
 484 and then extrapolate the data to a more exhaustive experiment, the tool could benefit from a completely

485 simulated mode, where tasks are generated dynamically using an average runtime of tasks input by the
486 user (and adapted with various factors to better represent the variability in runtime of a given experiment
487 and the execution on a distributed framework).

488 **6. Conclusions**

489 The developed framework allows speeding up hyperparameter optimization for medical image
490 classification significantly and easily (both for grid search and random sampling). The distributed nature
491 of the execution environment is leveraged for reducing the search space and gaining further wall-time.
492 The simulation tool allows estimating the runtime and results of medical texture analysis experiments
493 under various conditions, as well as extracting information such as a measure of the time–performance
494 trade–off of varying the cancellation margins. These tools can be used in a large variety of tasks that
495 include both image analysis and machine learning aspects. The system using Hadoop is relatively easy to
496 set up and we expect that many groups can make such optimizations in a much faster way using the results
497 of this article. Indeed, the dramatic reduction in runtime using only a local computing infrastructure can
498 enable the execution of experiments at a scale that may have been dismissed previously, ensuring to
499 get the best–possible results in the optimization of classification or similar tasks in a very reasonable
500 amount of time. The simulation environment can also help analyze performance and cost trade–offs
501 when optimizing parameters and potentially using cloud environments, allowing to give cost estimates.

502 The framework was evaluated with machine learning algorithms with a small number of
503 hyperparameters (i.e., two for SVMs and three for RFs). In future work, the framework is planned to
504 be tested with other datasets and more classifiers in order to validate its flexibility, potentially also with
505 approaches such as deep learning that can use several million hyperparameters and usually rely on GPU
506 computing [29], often supported by cloud providers as well. It is also planned to run comparative and
507 larger–scale experiments on a cloud–computing platform instead of using the local Hadoop infrastructure
508 to compare the influence of a mixed environment on runtime, as this can depend much more on
509 the available bandwidth. More advanced task cancellation criteria could also be implemented (e.g.
510 bandit–based method) to allow for more fine–grained control over the tasks to keep. Moreover, adding
511 more sophisticated parameter search strategies to the framework, such as Bayesian optimization or
512 gradient descent, could help improve the system even further, even though it will increase the complexity.

513 **Acknowledgments**

514 This work was supported by the Swiss National Science Foundation under grant PZ00P2_154891.

515 **Author Contributions**

516 Roger Schaer implemented the tools for running the grid parameter search in parallel, ran the
517 experiments, measured the results and wrote large parts of the article. Henning Müller provided ideas for
518 running the system and setting up the hardware infrastructure. Adrien Depeursinge provided the tested
519 image analysis and machine learning scenario and optimized the tools.

520 **Conflicts of Interest**

521 “The authors declare no conflict of interest”.

522 References

- 523 1. Kim, J. Iterated Grid Search Algorithm on Unimodal Criteria. PhD thesis, Virginia Polytechnic
524 Institute and State University, 1997.
- 525 2. Bergstra, J.; Bengio, Y. Random Search for Hyper-parameter Optimization. *J. Mach. Learn.*
526 *Res.* **2012**, *13*, 281–305.
- 527 3. Markonis, D.; Schaer, R.; Eggel, I.; Müller, H.; Depeursinge, A. Using MapReduce for
528 Large-scale Medical Image Analysis **2015**. [[arXiv:cs.DC/arXiv:1510.06937](https://arxiv.org/abs/cs/1510.06937)].
- 529 4. Owen, S.; Anil, R.; Dunning, T.; Friedman, E. *Mahout in Action*; Manning Publications Co.:
530 Greenwich, CT, USA, 2011.
- 531 5. Luo, G. MLBCD: a machine learning tool for big clinical data. *Health Information Science and*
532 *Systems* **2015**, *3*, 3.
- 533 6. Hutter, F.; Hoos, H.H.; Leyton-Brown, K. Sequential Model-based Optimization for General
534 Algorithm Configuration. Proceedings of the 5th International Conference on Learning and
535 Intelligent Optimization; Springer-Verlag: Berlin, Heidelberg, 2011; LION’05, pp. 507–523.
- 536 7. Friedrichs, F.; Igel, C. Evolutionary tuning of multiple SVM parameters. *Neurocomputing* **2005**,
537 *64*, 107–117.
- 538 8. Thornton, C.; Hutter, F.; Hoos, H.H.; Leyton-Brown, K. Auto-WEKA: Combined Selection
539 and Hyperparameter Optimization of Classification Algorithms. Proceedings of the 19th ACM
540 SIGKDD International Conference on Knowledge Discovery and Data Mining; ACM: New York,
541 NY, USA, 2013; KDD ’13, pp. 847–855.
- 542 9. Snoek, J.; Larochelle, H.; Adams, R.P. Practical Bayesian Optimization of Machine Learning
543 Algorithms. In *Advances in Neural Information Processing Systems 25*; Pereira, F.; Burges,
544 C.J.C.; Bottou, L.; Weinberger, K.Q., Eds.; Curran Associates, Inc., 2012; pp. 2951–2959.
- 545 10. Chapelle, O.; Vapnik, V.; Bousquet, O.; Mukherjee, S. Choosing Multiple Parameters for Support
546 Vector Machines. *Machine Learning* **2002**, *46*, 131–159.
- 547 11. Bergstra, J.S.; Bardenet, R.; Bengio, Y.; Kǎl’gl, B. Algorithms for Hyper-Parameter
548 Optimization. In *Advances in Neural Information Processing Systems 24*; Shawe-Taylor, J.;
549 Zemel, R.S.; Bartlett, P.L.; Pereira, F.; Weinberger, K.Q., Eds.; Curran Associates, Inc., 2011;
550 pp. 2546–2554.
- 551 12. Gorissen, D.; Couckuyt, I.; Demeester, P.; Dhaene, T.; Crombecq, K. A Surrogate Modeling
552 and Adaptive Sampling Toolbox for Computer Based Design. *J. Mach. Learn. Res.* **2010**,
553 *11*, 2051–2055.
- 554 13. Bergstra, J.; Komer, B.; Eliasmith, C.; Yamins, D.; Cox, D.D. Hyperopt: a Python library for
555 model selection and hyperparameter optimization. *Computational Science & Discovery* **2015**,
556 *8*, 014008.
- 557 14. Sparks, E.R.; Talwalkar, A.; Haas, D.; Franklin, M.J.; Jordan, M.I.; Kraska, T. Automating
558 Model Search for Large Scale Machine Learning. Proceedings of the Sixth ACM Symposium on
559 Cloud Computing; ACM: New York, NY, USA, 2015; SoCC ’15, pp. 368–380.

- 560 15. Liu, K.; Zhao, Q. Distributed Learning in Multi-Armed Bandit With Multiple Players. *IEEE*
561 *Transactions on Signal Processing* **2010**, *58*, 5667–5681.
- 562 16. Depeursinge, A.; Vargas, A.; Platon, A.; Geissbuhler, A.; Poletti, P.A.; Müller, H. Building a
563 Reference Multimedia Database for Interstitial Lung Diseases. *Computerized Medical Imaging*
564 *and Graphics* **2012**, *36*, 227–238.
- 565 17. Depeursinge, A.; Foncubierta-Rodríguez, A.; Van De Ville, D.; Müller, H. Multiscale Lung
566 Texture Signature Learning Using The Riesz Transform. *Medical Image Computing and*
567 *Computer-Assisted Intervention MICCAI 2012*. Springer Berlin / Heidelberg, 2012, Vol. 7512,
568 *Lecture Notes in Computer Science*, pp. 517–524.
- 569 18. Depeursinge, A.; Foncubierta-Rodríguez, A.; Van De Ville, D.; Müller, H. Rotation-covariant
570 texture learning using steerable Riesz wavelets. *IEEE Transactions on Image Processing* **2014**,
571 *23*, 898–908.
- 572 19. Dean, J.; Ghemawat, S. MapReduce: simplified data processing on large clusters. Proceedings
573 of the 6th conference on Symposium on Operating Systems Design & Implementation - Volume
574 6; USENIX Association: Berkeley, CA, USA, 2004; OSDI'04, pp. 10–10.
- 575 20. Shvachko, K.; Kuang, H.; Radia, S.; Chansler, R. The Hadoop Distributed File System.
576 Proceedings of the 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies
577 (MSST); IEEE Computer Society: Washington, DC, USA, 2010; MSST '10, pp. 1–10.
- 578 21. Vavilapalli, V.K.; Murthy, A.C.; Douglas, C.; Agarwal, S.; Konar, M.; Evans, R.; Graves,
579 T.; Lowe, J.; Shah, H.; Seth, S.; Saha, B.; Curino, C.; O'Malley, O.; Radia, S.; Reed, B.;
580 Baldeschwieler, E. Apache Hadoop YARN: Yet Another Resource Negotiator. Proceedings
581 of the 4th Annual Symposium on Cloud Computing; ACM: New York, NY, USA, 2013; SOCC
582 '13, pp. 5:1–5:16.
- 583 22. Hunt, P.; Konar, M.; Junqueira, F.P.; Reed, B. ZooKeeper: Wait-free Coordination for
584 Internet-scale Systems. Proceedings of the 2010 USENIX Conference on USENIX Annual
585 Technical Conference; USENIX Association: Berkeley, CA, USA, 2010; USENIXATC'10, pp.
586 11–11.
- 587 23. Hall, M.; Frank, E.; Holmes, G.; Pfahringer, B.; Reutemann, P.; Witten, I.H. The WEKA Data
588 Mining Software: An Update. *SIGKDD Explor. Newsl.* **2009**, *11*, 10–18.
- 589 24. Li, S.; Kwok, J.T.; Zhu, H.; Wang, Y. Texture classification using the support vector machines.
590 *Pattern Recognition* **2003**, *36*, 2883–2893.
- 591 25. Depeursinge, A.; Iavindrasana, J.; Hidki, A.; Cohen, G.; Geissbuhler, A.; Platon, A.; Poletti, P.A.;
592 Müller, H. Comparative Performance Analysis of State-of-the-Art Classification Algorithms
593 Applied to Lung Tissue Categorization. *Journal of Digital Imaging* **2010**, *23*, 18–30.
- 594 26. Vapnik, V.N. *The Nature of Statistical Learning Theory*; Springer: New York, 1995.
- 595 27. Breiman, L. Random Forests. *Machine Learning* **2001**, *45*, 5–32.
- 596 28. Quinlan, R.J. Induction of decision trees. *Machine Learning* **1986**, *1*, 81–106.
- 597 29. Krizhevsky, A.; Sutskever, I.; Hinton, G.E. ImageNet Classification with Deep Convolutional
598 Neural Networks. In *Advances in Neural Information Processing Systems 25*; Pereira, F.; Burges,
599 C.J.C.; Bottou, L.; Weinberger, K.Q., Eds.; Curran Associates, Inc., 2012; pp. 1097–1105.

600 © May 6, 2016 by the authors; submitted to *J. Imaging* for possible open access
601 publication under the terms and conditions of the Creative Commons Attribution license
602 <http://creativecommons.org/licenses/by/4.0/>.