

Communication

An Evolutionary Optimizer of *libsvm* Models

Dragos Horvath ^{1,*}, J.B. Brown ², Gilles Marcou ¹ and Alexandre Varnek ¹

¹ Laboratoire de Chémoinformatique, UMR 7140 CNRS – Univ. Strasbourg, 1, rue Blaise Pascal, 6700 Strasbourg, France

² Kyoto University Graduate School of Pharmaceutical Sciences, Department of Systems Bioscience for Drug Discovery, 606-8501, Kyoto, Japan

* Author to whom correspondence should be addressed; dhorvath@unistra.fr

Version August 7, 2014 submitted to *Challenges*. Typeset by \LaTeX using class file *mdpi.cls*

Abstract: This User Guide describes the rationale behind, and the *modus operandi* of a Unix script-driven package for evolutionary searching of optimal *libsvm* operational parameters, leading to Support Vector Machine models of maximal predictive power and robustness. Unlike common *libsvm* parameterizing engines, the current distributions includes the key choice of best-suited sets of attributes/descriptors, next to the classical *libsvm* operational parameters (kernel choice, cost, *etc.*), allowing a unified search in an enlarged problem space. It relies on an aggressive, repeated cross-validation scheme to ensure a rigorous assessment of model quality. Primarily designed for chemoinformatics applications, it also supports the inclusion of decoy instances – for which the explained property (bioactivity) is, strictly speaking, unknown, but presumably "inactive". The package supports four different parallel deployment schemes of the Genetic Algorithm-driven exploration of the operational parameter space, for workstations and computer clusters.

Keywords: chemoinformatics; QSAR; machine learning; *libsvm* parameter optimization

1. Introduction

Support Vector Machines (SVM), such as the very popular *libsvm* toolkit [1], are a method of choice for machine learning of complex, non-linear patterns of dependence of an explained variable – here termed the "property" P – and a set (vector) of attributes (descriptors \vec{D}) thought to be determinants of the current P value of the instance/object they characterize. The ultimate goal of machine learning

19 is to unveil the relationship between the property of an object and its descriptors: $P = P(\vec{D})$, where
20 this relationship may ideally represent a causal relationship between the attributes of the objects that
21 are responsible for its property, or, at least, a statistically validated covariance between attributes and
22 property. Following good scientific practice (Occam's razor), this relationship must take the simplest
23 form that explains so-far available experimental observations. Complexifying the relationship is only
24 justified if this leads to a significant increase of its predictive/extrapolation ability. Support Vector
25 Machines may describe relationships of arbitrary complexity. However, controlling the underlying
26 complexity of an SVM model may be only indirectly achieved, by fine tuning its operational parameters.
27 Unfortunately, this is a rather counterintuitive task (one cannot easily "guess" the optimal order of
28 magnitude for some of these parameters, not the mention the actual values), and cannot be addressed by
29 systematic scanning of all potentially valid combinations of control parameters. Furthermore, in various
30 domains, such as mining for Quantitative (Molecular) Structure – Property Relationships (QSPR), for
31 which this package has been primarily designed, there are many (hundreds) of *a priori* equally appealing,
32 alternative choices for the molecular descriptor set to be used as the input \vec{D} . There are many ways in
33 which the molecular structure can be encoded under the form of numeric descriptors. While some
34 problem-based knowledge may orient the user towards the most likely useful class of descriptors to
35 use in order to model a given molecular property P, tens to hundreds of empirical ways to capture the
36 needed molecular information in a bit-, integer- or real-number vector remain, and it cannot be known
37 beforehand which of these alternative "Descriptor Spaces" (DS) would allow for the most robust learning
38 of $P = P(\vec{D})$.

39 Moreover, vectors \vec{D} in these various DS may be of widely varying dimension, sparseness and
40 magnitude, meaning that some of the SVM parameters (in particular the γ coefficient controlling
41 Gaussian or Sigmoid kernel functions) will radically differ by many orders of magnitude with respect
42 to the employed DS. Since γ multiplies either a squared Euclidean distance, or a dot product value
43 between two vectors \vec{D} and \vec{d} in the considered DS in order to return the argument of an exponential
44 function, fitting will not focus on γ *per se*, but on a preliminary "gamma factor" γ_{gamma} defined such that
45 $\gamma = \gamma_{\text{gamma}} / \langle |\vec{D} - \vec{d}|^2 \rangle$ or, respectively $\gamma = \gamma_{\text{gamma}} / \langle \vec{D} \cdot \vec{d} \rangle$. Above, the denominators represent
46 means, over training set instance pairs, of Euclidean distance and dot product values, as required. In this
47 way, γ is certain to adopt a reasonable order of magnitude if the dimensionless γ_{gamma} parameters is
48 allowed to take values within (0,10).

49 The goal of the herein presented software distribution is to provide a framework for simultaneous
50 "meta"-optimization of relevant operational SVM parameters, including – and foremost – the choice of
51 best suited descriptor spaces, out of a list of predefined choices. This assumes simultaneous fitting
52 of categorical degrees of freedom (descriptor choice) and "classical" numerical SVM parameters,
53 where the acceptable orders of magnitudes for some of the latter directly depend on the former DS
54 choice. Such an optimization challenge – unfeasible by systematic search – clearly requires stochastic,
55 nature-inspired heuristics. A simple Genetic Algorithm (GA) has been implemented here. Chromosomes
56 are a concatenation (vector) of the currently used operational parameter values. To each chromosome,
57 a fitness score is associated, reporting how successful SVM model building was when the respective
58 choice of parameters was employed. Fitter setups are allowed to preferentially generate "offspring",
59 by cross-overs and mutations, and generate new potentially interesting setup configurations. Since this

60 is perfectly able to deal with both categorical and continuous degrees of freedom, additional toggle
 61 variables were considered: descriptor scaling (if set, initial raw descriptors of arbitrary magnitudes will
 62 first be linearly rescaled between zero and one) and descriptor pruning (if set, the initial raw descriptor
 63 space will undergo a dimensionality reduction, by removing one of pairwise correlated descriptor
 64 columns, respectively).

65 Model optimization is so-far supported for both ϵ Support Vector Regression (SVR) and Support
 66 Vector Classification (SVC) problems. The definition of SVR and SVC chromosome slightly differs,
 67 since the former also includes the ϵ tube width, within which regression errors are ignored. The nature
 68 of the problem may be indicated by command line parameter.

69 The fitness (objective) function (SVM model quality score) to be optimized by picking the
 70 most judicious parameter combination should accurately reflect the ability of the resulting model to
 71 extrapolate/predict instances not encountered during the training stage, not its ability to (over)fit training
 72 data. In this sense, "classical" leave-1/ N -out cross-validation (XV) already implemented in libsvm was
 73 considered too lenient, because its outcome may be strongly biased by the peculiar order of instances in
 74 the training set, dictating the subsets of instances that are simultaneously being left out. A M -fold
 75 repeated leave-1/ N -out XV scenario, where each M -fold repeat proceeds on a randomly reordered
 76 training set item list, has been retained here. As a consequence, $M \times N$ individual SVM models are by
 77 default built at any given operational parameter choice as defined by the chromosome. However, if any of
 78 these $M \times N$ attempts fails at fitting stage (quality of fit below a user-defined threshold), then the current
 79 parameterization is immediately discarded as unfit, in order to save computer time. Otherwise, the
 80 quality score used as fitness function of the GA not only reflects the average cross-validation propensity
 81 of models at given operational setup, but also accounts for its robustness with respect to training set
 82 item ordering. At equal average XV propensity over the M trials, models that cross-validate roughly
 83 equally well irrespective of the order to of training set items are to be preferred over models which
 84 achieve sometimes very strong, other times quite weak XV success depending on how left-out items
 85 were grouped together.

86 For SVR models, the cross-validated determination coefficient Q^2 is used to define fitness, as follows:

- 87 1. For each XV attempt between 1... M , consider the N models generated by permuting the left-out
 88 $1/N$ of the training set.
- 89 2. For each training instance i , retrieve the predicted property $\hat{P}_M(i)$ returned by the one specific
 90 model among the above-mentioned N which had not included i in its training set.
- 91 3. Calculate the cross-validated Root-Mean-Squared Error $RMSE_M$ of stage M , based on the
 92 deviations of $\hat{P}_M(i)$ from the experimental property of instance i . Convert this into the
 93 determination coefficient $Q_M^2 = 1 - RMSE_M^2/\sigma^2(P)$, by reporting this error to the intrinsic
 94 variance $\sigma^2(P)$ of property P within the training set.
- 95 4. Calculate the mean $\langle Q_M^2 \rangle$ and the standard deviation σ of Q_M^2 values over the M attempts.
 96 Then, fitness is defined as $\langle Q_M^2 \rangle - \kappa \times \sigma(Q_M^2)$, with a user-defined κ (2, by default).

97 In SVC, the above procedure is applied to the Balanced Accuracy (BA) of classification, defined as the
 98 mean, over all classes, of the fractions of correctly predicted members of class C out of the total number

99 of instances in class C. This is a more rigorous criterion than the default fraction of well-classified
100 instances reported by libsvm (sum of well-classified within each class/training set size), in particular
101 with highly unbalanced sets, as often the case in chemoinformatics.

102 Eventually, the set of $M \times N$ individual SVM models may serve as final consensus model for external
103 predictions if the fitness score exceeds a user-defined threshold and if, of course, external descriptor sets
104 are being provided. Furthermore, if these external sets also feature activity data, external prediction
105 challenge quality scores are calculated on-the-fly. However, the latter are not used in the Darwinian
106 selection process of the optimal parameter chromosome, but merely generated for further exploitation
107 by the user. External "on-the-fly" validation also allows to eventually delete the $M \times N$ temporarily
108 generated models, avoiding the potential storage problems that would be quick to emerge if several
109 tens of bulky SVM model files were to be kept for each of the thousands of attempted parameterization
110 schemes. The user is given the option to eventually explicitly regenerate the model files for one or more
111 of the most successful parameterization schemes "discovered" by the Darwinian process, and use them
112 outside of the context of this package.

113 Both the expansion of parameter space to include "meta"-variables such as the chosen DS and the
114 pretreatment protocol of descriptors and the aim for a robust model quality estimated based on an
115 M-fold more expensive approach than classical XV may render the mining for the optimal model
116 quite effort-consuming. Therefore, the current distribution supports both a multi-CPU workstation
117 version for "easy" problems (hundreds of training items within tens of candidate DS), but also parallel
118 deployment-ready pilots using Torque, Slurm or LSF (Load Share Facility) on clusters.

119 2. Installation and Use

120 This distribution includes a set of tcsh scripts, relying on various subscripts (tcsh, awk and perl) and
121 executables in order to perform specific tasks. This User guide assumes basic knowledge of Linux – it
122 will not go into detailed explanations concerning fundamentals such as environment, paths, *etc.*

123 2.1. Installation and Environment Setup

124 2.1.1. Prerequisites

125 As a prerequisite, check whether `tcsh`, `awk` and `perl` are already installed on your system. Also,
126 the `at` job scheduler should be enabled, both on the local workstations and on cluster nodes, as needed.
127 In principle, our tools should not be tributary to specific versions of these prerequisite standard packages,
128 and such dependence was never observed so far. Note, however, that `awk` (unfortunately) displays a
129 language-dependent behavior, and may be coaxed to accept a comma as decimal separator instead of
130 the dot, which would produce weird results (as `awk` always does something, and never complains).
131 Therefore, we explicitly recommend to check the language setup of your machine – if in doubt, set your
132 `LANG` environment variable to "en_US". Also, if you plan to machine-learn from DOS/Windows files
133 and forgot to convert them to Unix-style, expect chaotic behavior. Install the `dos2unix` utility now, so
134 you won't have to come back to this later.

135 2.1.2. Installation

136 Download the `.tar.gz` archive and extract files in a dedicated directory. It does not include the
 137 *libsvm* package, which you must install separately. Two environment variables must be added to your
 138 configuration: `GACONF`, being assigned the absolute path to the dedicated directory, and `LIBSVM`, set
 139 to the absolute path of the residence directory of *libsvm* tools. An example (for `tcsh` shells) is given in
 140 the file **EnvVars.csh** of the distribution, to be edited and added to your `.cshrc` configuration file. Adapt
 141 as needed (use `export` instead of `setenv`, in your `.profile` configuration file) if your default shell is other
 142 than `(t) csh`.

143 Specific executables shipped with this distribution are compiled for `x86_64` Linux machines. They
 144 are all regrouped in the subdirectory `$GACONF/x86_64`, since it is assumed that your machine type, as
 145 encoded by the environment variable `$MACHTYPE`, is `"x86_64"`. If this is not your case, you will have
 146 to create your own executable subdirectory `$GACONF/$MACHTYPE` and generate the corresponding
 147 executables in that directory. To this purpose, source codes are provided for all Fortran utilities. They
 148 can be straightforwardly recompiled:

```
149
150 foreach f ($GACONF/*.f)
151 g77 -O3 -o $GACONF/$MACHTYPE/$f:r $f
152 end
```

153
 154 – please adapt to your favorite Unix shell syntax. Note – `gfortran`, `f77` or other Fortran
 155 compilers may do the job as well: those snippets of code are fairly robust, using standard commands.

156 However, *simScorer*, the executable used to calculate mean Euclidean distance and dot product values
 157 of training set descriptor vectors is a full-blown, strategic virtual screening tool of our laboratory. Its
 158 FreePascal source code cannot be released. However, we might perhaps be able to generate, upon
 159 request, the kind of executable you require.

160 2.1.3. Deployment-Specific Technical Parameters

161 As hinted in the Introduction, this distribution supports a parallel deployment of the fitness scoring of
 162 operational parameter combinations encoded by the chromosomes. Four scenarios are supported:

163 1. **local:** several CPUs of a workstation are used in parallel, each hosting an independent machine
 164 learning/XV attempt (the "slave" job) with parameters encoded by the current chromosome, and
 165 returning, upon completion, the fitness score associated to that chromosome to the master "pilot"
 166 script overseeing the process. Master (pilot) and slave (evaluators of the parameterization scheme
 167 quality) scripts run on the same machine – however, since the pilot sleeps for most of the time,
 168 waiting for slaves to return their results, the number of slave jobs may equal the number of CPU
 169 cores available. As the number of active slave processes decreases, the pilot will start new ones,
 170 until the user-defined maximal number of parameterization attempts has been reached.

171 2. **torque:** slaves estimating the quality of a parameterization scheme are being submitted on
 172 available nodes of a cluster running the torque job scheduler. The pilot runs on the front-end

(again, without significantly draining computer power) and schedules a novel slave job as soon as the number of executing or waiting slave jobs has dropped below the user-defined limit, until the user-defined maximal number of parameterization attempts has been reached. One slave job is scheduled to use a single CPU core (meaning that, like in the local version, successive XV runs are executed sequentially). Slave jobs are being allotted explicit time lapses (walltime) by the user and those failing to complete within allotted time are lost.

3. **slurm**: unlike in torque-driven deployment, the pilot script running on the front-end of the slurm-operated batch system reserves entire multi-CPU nodes for the slave jobs, not individual CPUs. Since slave jobs were designed to run on a single CPU, the slurm pilot does not directly manage them, but deploys, on each reserved node, a slave manager mandated to start the actual slave jobs on each CPU node, then wait for their completion. When the last of the slaves finished and returned results to the front-end pilot, the manager completes as well, and frees the node. The front-end pilot observes this, and reschedules another session on a free node, until the user-defined maximal number of parameterization attempts is reached. A runtime limit per node (walltime) must be specified: when reached, the manager and all so-far incomplete slave jobs will be killed.
4. **bsub**: this batch system (officially known as **LSF**, but practically nicknamed **bsub** by the name of its key job submission command) is similar to slurm, in the sense that the front-end pilot is managing multi-CPU nodes, not individual cores. However, there is no mandatory time limit per node involved. The task of the front-end pilot is to start a user-defined number of "local" workstation pilot jobs on as many nodes, keep track of incoming results, potentially reschedule local pilots if some happened to terminate prematurely, and eventually stop the local pilots as soon as the user-defined maximal number of parameterization attempts has been reached. Local pilots will continue feeding new slave jobs to the available CPUs of the node, not freeing the current node until terminated by the front-end pilots (or by failure).

IMPORTANT! Cluster-based strategies may require some customization of the default parameters. These may be always modified at command-line startup of the pilot scripts, but permanent modifications can be hard-coded by editing the corresponding ***.oppars** files ("*****" meaning **torque**, **slurm**, **bsub**, respectively) in `$GACONF`. The files contain commented set instructions (`tcsh`-style, do *not* translate to your default shell syntax) and should be self-explanatory. **Some adjustments – account and queue specifications, when applicable – are mandatory.** Moreover, you might want to modify the `tmpdir` variable, which points to the local path on a node-own hard drive, on which temporary data may be written. Unless it crashed, the slave job will move relevant results back to your main work directory, and delete the rest of temporary data. Do not modify the `restart` parameter (explained further on).

2.2. *User Guide: Understanding and Using the libsvm Parameter Configurator* The `libsvm` parameter configurator is run by invoking, at command line, the pilot script dedicated to the specific deployment scheme supported by your hardware:

```
210
211 $GACONF/pilot_<scheme>.csh data_dir=<directory containing input
212 data> workdir=<intended location of results> [option=value...] >&
213 logfile.msg
```

214
215 where <scheme> is either of the above-mentioned *local*, *torque*, *slurm*, *bsub*. The
216 script does not attempt to check whether it is actually run on a machine supporting the envisaged
217 protocol – calling *pilot_slurm.csh* on a plain workstation will result in explicit errors, while calling
218 *pilot_local.csh* on the front-end of a slurm cluster will charge this machine with machine learning
219 jobs, much to the disarray of other users and the ire of the system manager.

220 2.2.1. Preparing the Input Data Directory

221 The minimalistic input of a machine learning procedure is a property-descriptor matrix of the training
222 set instances. However, one of the main goals of the present procedure is to select the best suited set of
223 descriptors, out of a set of predefined possibilities. Therefore, several property-descriptor matrices (one
224 per candidate descriptor space) must be provided. As a consequence, it was decided to regroup all the
225 required files into a directory, and pass the directory name to the script, which then will detect relevant
226 files therein – following naming conventions. An example of input directory, **D1-datadir**, is provided
227 with the distribution. It is good practice to keep a copy of the list of training set items in the data
228 directory, even though this will not be explicitly used. In **D1-datadir**, this is entitled **ref.smi_act_info**.
229 It is a list of SMILES strings (column #1) of 272 ligands for which experimental affinity constants (pK_i ,
230 column #2) with respect to the D1 dopamine receptor were extracted from the ChEMBL database [2],
231 forming the training set of a D1 affinity prediction model.

232 Several alternative schemes to encode the molecular information under numeric form were
233 considered, *i.e.* these molecules were encoded as position vector in several DS. These DS include
234 pharmacophore (PH)- and atom-symbol (SY) labeled sequences (seq) and circular fragment counts
235 (tree,aab), as well as fuzzy pharmacophore triplets [3,4]. For each DS, a descriptor file encoding
236 one molecule per line (in the order of the reference list) is provided: **FPT1.svm**, **seqPH37.svm**,
237 **seqSY37.svm**, **aabPH02.svm**, **treeSY03.svm**, **treePH03.svm** in .svm format. Naming of these
238 descriptor files should consist of an appropriate DS label (a name of letters, numbers and underscores –
239 *no spaces, no dots* – unequivocally reminding you how the descriptors were generated), followed by the
240 **.svm** extension. For example, if you have three equally trustworthy/relevant commercial pieces of soft,
241 each offering to calculate a different descriptor vector for the molecules, let each generate its descriptors
242 and store them as **Soft1.svm**, **Soft2.svm**, **Soft3.svm** in the data directory. Let us generically refer to
243 these files as **DS.svm**, each standing for a descriptor space option. These should be plain Unix text files,
244 following the SVM format convention, *i.e.* representing each instance (molecule) on a line, as:

```
245
246 SomeID 1:Value_Of_Descriptor_Element_1 2:Value_Of_Descriptor_Element_2
247 ... n:Value_Of_Descriptor_Element_n
```

248

249 where the values of vector element D_i refer, of course, to the current instance described by the current
250 line. The strength of this way to render the vector \vec{D} is that elements equal to zero may be omitted, so
251 that in terms of total line lengths the explicit citing of the element number ahead of the value is more
252 than compensated for by the omission of numerous zeros in very sparse vectors. If your software does
253 not support writing of .svm files, it is always possible to convert a classical tabular file to .svm format,
254 using:

```
255  
256 awk '{printf "%s", $1; for (i=2; i<=NF; i++) {if ($i!=0) printf  
257 "%d:%.3f", i-1, $i}; print ""}' tabular_file_with_first_ID_column.txt  
258 > DS_from_tabular.svm
```

259
260 where you may tune the output precision of descriptors (here %.3f). Note – there are no title/header
261 lines in .svm files. If the tabular file does contain one, make sure to add an NR>1 clause to the above
262 awk line.

263 By default, .svm files do not contain the modeled property as the first field on the line, as expected
264 from a property-descriptor matrix. Since several .svm files, each representing another DS – and
265 presumably generated with another piece of soft, that may or may not allow the injection of the modeled
266 property parameter into the first column – our software tools allow .svm files with first columns of
267 arbitrary content to be imported as such into the data directory. The training set property values need
268 to be specified separately, in a file having either the extension .SVMreg – for continuous properties to
269 be modeled by regression, or SVMclass, for classification challenges. You may have both types of
270 files, but only one file of a given type per data directory (their name is irrelevant, only the extension
271 matters). For example, in **D1-datadir**, the file **ref.SVMreg** is a single-column file of associated D1
272 pK_i affinity values (matching the order of molecules in descriptor files). The alternative **ref.SVMclass**
273 represents a two-class classification scheme, distinguishing "actives" of $pK_i > 7$ from "inactives". The
274 provided scripts will properly merge property data with corresponding descriptors from .svm files during
275 the process of cross-validated model building. Cross-validation is an automated protocol – the user needs
276 not to provide a split of the files into training and test sets: just upload the global descriptor files to the
277 data directory.

278 In addition, our package allows on-the-fly external predictions and validation of the models. This is
279 useful in as far as, during the search stage of the optimal parameters, models generated according to a
280 current parameter set are evaluated, then discarded, in order not to fill the disks with bulky model files that
281 seemed promising at the very beginning of the evolutionary process, but were long since outperformed
282 by later generations of parameter choices. Do not panic: the model files corresponding to user-picked
283 parameter configurations may always be recreated and kept for off-line use. However, if descriptor files
284 (or property-descriptor matrices) of external test sets are added to the data directory, models of fitness
285 exceeding a user-defined threshold will be challenged to predict properties for the external instances
286 before deletion. The predictions, of course, will be kept and reported to the user, but never used in
287 parameter scheme selection – the latter being strictly piloted by the cross-validated fitness score. Since
288 the models may be based on either of the initially provided **DS.svm**, all the corresponding descriptor
289 sets must also be provided for each of the considered external test sets. In order to avoid confusion with

290 training files, external set descriptors must be labeled as **ExternalSetName.DS.psvm**. Even if only one
291 external set is considered, the three-part dot-separated syntax must be followed (therefore, no dots within
292 ExternalSetName or DS, please). For example, check the external prediction files in **D1-datadir**. They
293 correspond to an external set of dopamine D5-receptor inhibitors, encoded by the same descriptor types
294 chosen for training: **D5.FPT1.psvm**, **D5.seqPH37.psvm**, ... Different termination notwithstanding,
295 **.psvm** files are in the same **.svm** format. Whatever the first column of these files may contain, our
296 tool will attempt to establish a correlation between the predicted property and this first column (forcibly
297 interpreted as numeric data). If the user actually provided the experimental property values, then external
298 validation statistics are automatically generated. Otherwise, if these experimental values are not known
299 – this is a genuine prediction exercise – then the user should focus on the returned prediction values.
300 Senseless external validation statistics scores may be reported if your first column may be interpreted as
301 a number – just ignore. In the sample data provided here, the first column reports pK_i affinity constants
302 for the D5 receptor, whereas the predicted property will be the D1 affinity. This is thus not an actual
303 external validation challenge, but an attempt to generate computed D1 affinities for D5 ligands.

304 To wrap up, the input data directory must/may contain the following files:

- 305 1. A **list** of training items (optional)
- 306 2. A one-column **property** file with extensions **.SVMreg** (regression modeling of continuous
307 variables) or **SVMclass** (classification modeling) respectively (**compulsory**)
- 308 3. **Descriptor** files, in **.svm** format, one per considered DS, numerically encoding one training
309 instance per line, ordered like in the property file (**at least one DS.svm required**). The contents
310 of their first column is irrelevant, as the property column will be inserted instead.
- 311 4. Optional **external** prediction files – for each considered external set, all the associated descriptors
312 must be provided as **ExternalSetName.DS.psvm** files. The user has the option of inserting
313 experimental properties to be automatically correlated with predictions in the first column of these
314 **.psvm** files

315 2.2.2. Adding Decoys – The Decoy Directory

316 In chemoinformatics, it is sometimes good practice to increase the diversity of a con-generic training
317 set (based on a common scaffold) by adding very different molecules. Otherwise, machine learning may
318 fail to see the point that the common scaffold is essential for activity and learn the "antipharmacophore"
319 – features that render some of these scaffold-based analogues *less* active than others. Therefore, when
320 confronted to external predictions outside of the scaffold-based families, such models will predict all
321 compounds *not* containing this antipharmacophore (including cosmic vacuum, matching the above
322 condition) to be active. It makes sense to teach models that outside the scaffold-based family activity will
323 be lost. This may be achieved by adding *presumed* diverse inactives to the training set. However, there is
324 no experimental (in)activity measure for these – empirically, a value smaller than the measured activity
325 of the less active genuine training set compound is considered. This assumption may, however, distort
326 model quality assessment. Therefore, in this approach we propose an optimal compromise scenario in
327 which, if so desired by the user, the training set may be augmented by a number of decoy instances

328 equal to half of the actual training set size. At each XV attempt, after training set reshuffling, a (every
329 time different) random subset of decoys from a user-defined **decoy repository** will be added. The SVM
330 model fitting will account for these decoys, labeled as inactives. However, the model quality assessment
331 will ignore the decoy compounds, for which no robust experimental property is known: it will focus
332 on the actual training compounds only. Therefore, the user should *not* manually complete the training
333 set compounds in the data directory with decoy molecules. An extra directory – to be communicated to
334 the pilot script using the `decoy_dir=<path of decoy repository>` command line option –
335 should be set up, containing **DS.svm** files of an arbitrary large number of decoy compounds. Again, all
336 the DS candidates added to the data directory must be represented in the decoy directory. The software
337 will therefrom extract random subsets of size comparable to the actual training set, automatically inject
338 a low activity value into the first column, and join these to the training instances, creating a momentarily
339 expanded training set.

340 In case of doubt concerning the utility of decoys in the learning process, the user is encouraged to run
341 two simulations in parallel – one employing decoys, the other not (this is the default behavior unless
342 `decoy_dir=<path to decoy repository>` is specified). Unless decoy addition strongly
343 downgrades model fitness scores, models employing decoys should be preferred because they have an
344 intrinsically larger applicability domain, being confronted with a much larger chemical subspace at the
345 fitting stage.

346 2.2.3. Command-line Launch of the libsvm Parameter Configurator: the Options and their Meaning

347 The following is an overview of the most useful options that can be passed to the pilot script, followed
348 by associated explanations of underlying processes, if relevant:

- 349 1. `workdir=<intended location of results>` is always mandatory, for both new
350 simulations (for which the working directory must not already exist, but will be created) and
351 simulation restarts, in which the working directory already exists, containing preprocessed input
352 files and so-far generated results.
- 353 2. `cont=yes` must be specified in order to restart a simulation, based on an existing working
354 directory – in which data preprocessing is completed. Unless this flag is set (default is new run),
355 specifying an existing working directory will result in an error.
- 356 3. `data_dir=<directory containing input data>` is mandatory for each new
357 simulation, for it refers to the repository of basic input information, as described above. Unless
358 `cont=yes`, failure to specify the input repository will result in failure.
- 359 4. `decoy_dir=<path of decoy repository>` as described in the previous subsection is
360 mandatory only for new starts (if `cont=yes`, preprocessed copies of decoy descriptors are
361 assumed to already reside in the working directory).
- 362 5. `mode=SVMclass` toggles the pilot to run in classification mode. Default is `mode=SVMreg`, *i.e.*
363 *libsvm* ϵ -regression. Note: the extension of the property data file in the input data directory must
364 strictly match the `mode` option value (it is advised *not* to mistype "SVMclass").

365 6. `prune=yes` is an option concerning the descriptor preprocessing stage, which will be described
366 in this paragraph. This preprocessing is the main job of `$GACONF/svmPreProc.pl`, operating
367 within `$GACONF/prepWorkDir.csh`, called by the pilot script. It may perform various
368 operations on the brute **DS.svm**. By default, training set descriptors from the input directory
369 are scanned for columns of near-constant or constant values. These are discarded if the standard
370 deviation of vector element i over all training set instances is lower than 2% of the interval
371 $maxD_i - minD_i$ (near-constant), or if $maxD_i = minD_i$ (constant). Then, training set descriptors
372 are Min/Max-scaled (for each descriptor column i , $maxD_i$ is mapped to 1.0, while $minD_i$ maps
373 to 0.). This is exactly what the *libsvm*-own `svm-scale` would do – now it is achieved on-the-fly,
374 within the larger preprocessing scheme envisage here. Training set descriptors are the ones used
375 to fix the reference minimum and maximum, further used to scale all other terms (external set
376 `.psvm` files and/or decoy `.svm` files). These reference extremes are stored in `.pri` files in the
377 working directory, for further use. In this process, descriptor elements ("columns") are sorted with
378 respect to their standard deviations and renumbered. If, furthermore, the `prune=yes` option has
379 been invoked, a (potentially time-consuming) search for pairwise correlated descriptor columns
380 (at $R^2 > 0.7$) will be performed, and one member of each concerned pair will be discarded.
381 However, it is not *a priori* granted that either Min/Max scaling or pruning of correlated columns
382 will automatically lead to better models. Therefore, scaling and pruning are considered as degrees
383 of freedom of the GA – the final model fitness is the one to decide whether scaled, pruned, scaled
384 and pruned or plain original descriptor values are the best choice to feed into the machine learning
385 process. At preprocessing stage, these four different versions (scaled, pruned, scaled and pruned
386 or plain original) are generated, in the working directory, for each of the `.svm` and `.psvm` files in
387 input or decoy folders. If pruning is chosen, but for a given DS the number of cross-correlated
388 terms is very low (less than 15% of columns could be discarded in the process), then pruning is
389 tacitly ignored – the two "pruned" versions are deleted because they would likely behave similarly
390 to their parent files. For each of the processed descriptor files – now featuring property values in
391 the first column, as required by `svm-train` – the mean values of descriptor vector dot products,
392 respectively Euclidean distances are calculated and stored in the working directory (file extensions
393 **.EDprops**). For large training sets (> 500 instances), these means are not taken over all the
394 $N(N - 1)/2$ pairs of descriptor vectors, but are calculated on hand of a random subset of 500
395 instances only.

396 **IMPORTANT!** The actual descriptor files used for model building – and expected as input for
397 model prediction – are, due to preprocessing, significantly different from the original **DS.svm**
398 generated by yourself in the input folder. Therefore, any attempt to use generated models in a
399 stand-alone context, for property prediction, must take into account that any input descriptors must
400 first be reformatted in the same way in which an external test `.psvm` file is being preprocessed.
401 Calling the final model on brute descriptor files will only produce noise. All the information
402 required for preprocessing is contained in `.pri` files. Suppose, for example, that Darwinian
403 evolution in *libsvm* parameter space showed that the so-far best modeling strategy is to use the
404 pruned and scaled version of `DS.svm`. In that case, in order to predict properties of external
405 instances described in **ExternalDS.svm**, one would first need to call:

406
407 `$GACONF/svmPreProc.pl ExternalDS.svm selfile=DS_pruned.pri`
408 `scale=yes output=ReadyToUseExternalDS.svm`
409

410 where the required **DS_pruned.pri** can be found in the working directory. If original (not scaled)
411 descriptors are to be used, drop the `scale=yes` directive. If pruning was not envisaged (or shown
412 not to be a good idea, according to the undergone evolutionary process), use the **DS.pri** selection
413 file instead. The output `.svm` is now ready to be used for prediction by the generated models.

414 7. `wait=yes` instructs the pilot script to generate the new working directory and preprocess input
415 data, as described above, and then stop rather than launching the GA. The GA simulation may be
416 launched later, using the `cont=yes` option.

417 8. `maxconfigs=<maximal number of parameter configurations to be`
418 `explored>` – by default, it is set to 3000. It defines the amount of effort to be invested
419 in this search of the optimal parameter set, and should take into account the number of considered
420 descriptor spaces impacting on the total volume of searchable problem space. However, not being
421 able to guess the correct `maxconfigs` values is not a big problem. If the initial estimate seems
422 to be too high, no need to wait for the seemingly endless job to complete: one may always trigger
423 a clean stop of the procedure by creating a file (even an empty one is fine) named **stop_now** in the
424 working directory. If the estimate was too low, one may always apply for more number crunching
425 by restarting the pilot with options `workdir=<already active working directory>`
426 `cont=yes maxconfigs=<more than before>`

427 9. `nnodes=<number of "nodes" on which to deploy simultaneous slave`
428 `jobs>` is a possibly misleading name for a context-dependent parameter. For the `slurm` and
429 `LSF` contexts, this refers indeed to the number of multi-CPU machines (nodes) required for
430 parallel parameterization quality assessments. Under the `torque` batch system and on local
431 workstations, it actually stands for the number of CPU cores dedicated to this task. Therefore, the
432 `nnodes` default value is context-dependent: 10 for `slurm` and `LSF/bsub`, 30 for `torque` and
433 equal to the actual number of available cores on local workstations.

434 10. `lo=<"leave-out" XV multiplicity>` is an integer parameter greater or equal to two,
435 encoding the number of folds into which to split the training set for XV. By default, `lo=3`, meaning
436 that 1/3 of the set is iteratively kept out for predictions by a model fitted on the remaining 2/3
437 of training items. Higher `lo` values make for easier XV, as a larger part of data is part of the
438 momentary training set, producing a model having "seen" more examples and therefore more
439 likely to be able to properly predict the few remaining 1/`lo` test examples. Furthermore, since `lo`
440 obviously gives the number of model fitting jobs needed per XV cycle, higher values translate to
441 proportionally higher CPU efforts. If your data sets are small, so that 2/3 of it would likely not
442 contain enough information to support predicting the left-out 1/3, you may increase `lo` up to 5. Do
443 not go past that limit – you are only deluding yourself by making XV too easy a challenge, and
444 consume more electricity, atop of that.

445 11. `ntrials`=<number of repeated XV attempts, based on randomized
446 regrouping of kept and left-out subsets> dictates how many times (12,
447 by default) the leave-1/lo-out procedure has to be repeated. It matches the formal parameter M
448 used in Introduction in order to define the model fitness score. The larger `ntrials`, the more
449 robust the fitness score (the better the guarantee that this was not some lucky XV accident due to a
450 peculiarly favorable regrouping of kept vs. left-out instances.) However, the more time-consuming
451 the simulation will be, as the total number of fitted local models during XV equals `ntrials` \times
452 `lo`. A possible escape from this dilemma of quick vs. rigorous XV is to perform a first run over a
453 large number (`maxconfigs`) of parameterization attempts, but using a low `ntrials` XV repeat
454 rate. Next, select the few tens to hundreds of best-performing parameter configurations visited so
455 far, and use them (see option `use_chromo` below) to rebuild and reestimate the corresponding
456 models, now at a high XV repeat rate.

457 12. `use_chromo`=<file of valid parameter configuration chromosomes> is
458 an option forcing the scripts to create and assess models at the parameter configurations from
459 the input file, rather than using the GA to propose novel parameter setup schemes. These valid
460 parameter configuration chromosomes have supposedly emerged during a previous simulation
461 – therefore, the `use_chromo` option implicitly assumes `cont=yes`, *i.e.* an existing working
462 directory with preprocessed descriptors. As will be detailed below, the GA-driven search of
463 valid parameter configurations creates, in the working directory, two result files: **done_so_far**,
464 reporting every so-far assessed parameter combination and the associated model fitness criteria,
465 and **best_pop**, a list of the most diverse setups among the so-far most successful ones.

466 There are many more options available, which will not be detailed here because their default values rarely
467 need to be tampered with. Geeks are encouraged to check them out, in comment-adorned files ending in
468 ***pars** provided in `$GACONFIG`. There are common parameters (**common.pars**), deployment-specific
469 parameters in ***.oppars** files, and model-specific (*i.e.* regressions-specific vs. classification-specific)
470 parameters **SVMreg.pars**, **SVMclass.pars**. The latter concern two model quality cutoffs

- 471 1. `minlevel` – only models better than this, in terms of fitness scores, will be submitted to external
472 prediction challenges.
- 473 2. `fit_no_go` represents the minimal performance at fitting stage, for every local model built
474 during the XV process. If (by default) a regression model fitting attempt fails to exceed a fit
475 R^2 value (or a classification model fails to discriminate at better BA), then hope to see this model
476 succeed in terms of XV is low. In order to save time, the current parameterization attempt is
477 aborted – the parameter choice is obviously wrong – and time is saved by reporting a fictitious,
478 very low fitness score associated to the current parameter set.

479 Default choices for regression models refer to correlation coefficients, and are straightforward to
480 interpret. However, thresholds for classification problems depend on the number of classes. Therefore,
481 defaults in **SVMclass.pars** correspond not to absolute balanced accuracy levels, but to fractions of
482 non-random BA value range (parameter value 1.0 maps to a BA cutoff value of 1, while parameter
483 value zero maps to the baseline BA value, equaling the reciprocal of the class number).

484 2.2.4. Defining the Parameter Phase Space

485 As already mentioned, the model parameter space to be explored features both categorical and
486 continuous variables. The former include DS and *libsvm* kernel type selection, the latter cover ϵ (for
487 SVMreg, *i.e.* ϵ -regression mode only), cost γ and *coeff0* values. In general, as well as in this
488 case, GAs typically treat continuous variables like a user precision-dependent series of discrete values,
489 within a user-defined range. The GA problem space is defined in two mode-dependent setup files in
490 `$GACONF: SVMreg.rng` and `SVMclass.rng`, respectively. They feature one line for each of the
491 considered problem space parameters, except for the possible choices of DS, which are not available
492 by default. After the preprocessing stage, when standardized `.svm` files were created in the working
493 directory, the script will make a list of all the available DS options, and write it out as the first line
494 of a local copy (within the working directory) of the `<mode>.rng` file. Then it will concatenate the
495 default `$GACONF/<mode>.rng` to the latter. This local copy is the one piloting the GA simulation,
496 and may be adjusted whenever the defaults from `$GACONF/<mode>.rng` seem inappropriate. To do
497 so, first invoke the pilot script with the data repository, a new working directory name, desired mode and
498 option `wait=yes`. This will create the working directory, uploading all files and generating the local
499 `<mode>.rng` file, then stop. Edit this local file, then re-invoke the pilot on the working directory, with
500 option `cont=yes`.

501 The syntax of `.rng` files is straightforward. Column one designs the name of the degree of freedom
502 defined on the current line. If this degree of freedom is categorical, the remaining fields on the
503 line enumerate the values it may take. For example, line #1 is a list of descriptor spaces (including
504 their optionally generated "pruned" versions) found in the working directory. The parameter "scale"
505 chooses whether those descriptors should be used in their original form, or after Min/Max scaling
506 (in clear, if scale is "orig", then `DS.orig.svm` will be used instead of `DS.scaled.svm`). The "kernel"
507 parameter picks the kernel type to be used with *libsvm*. Albeit a numeric code is used to define it
508 on the `svm-train` command line (0- linear, 1 - 3rd order polynomial, 2 - Radial Basis Function, 3-
509 Sigmoid), on the `.rng` file line, the options were prefixed by "k" in order to let the soft handle this
510 as a categoric option. The "ignore-remaining" keyword on the kernel line is a directive to the GA
511 algorithm to ignore the further options γ and *coeff0* if the linear kernel choice `k0` is selected. This
512 is required for population redundancy control: two chromosomes encoding identical choices for all
513 parameters except γ and *coeff0* do stand for genuinely different parameterization schemes, leading to
514 models of different quality – *unless* the kernel choice is set to linear, which is not sensitive to γ and
515 *coeff0* choices, leading to exactly the same modeling outcome. All `.rng` file lines having a numeric
516 entry in column #2 are by default considered to encode continuous (discretized) variables. Such lines
517 must obligatorily have 6 fields (besides #-marked comments): variable name, absolute minimal value,
518 preferred minimal value, preferred maximum, absolute maximum and, last, output format definition
519 implicitly controlling its precision. The distinction between absolute and preferred ranges is made in
520 order to allow the user to focus on a narrowest range assumed to contain the targeted optimal value, all
521 while leaving an open option for the exploration of larger spaces: upon random initialization or mutation
522 of the variable, in 80% of cases the value will be drawn within the "preferred" boundaries, in 20% of
523 cases within the "absolute" boundaries. The brute real value is then converted into a discrete option

524 according to the associated format definition in the last column – it is converted to a string representation
525 using `printf("<format representation>", value)`, thus rounded up to as many decimal
526 digits as mentioned in the decimal part of its format specifier. For example, the cost parameter spanning a
527 range of width 18, with a precision of 0.1 may globally adopt 180 different values. Modifying the default
528 "4.1" format specifier to "4.2" triggers a 10-fold increase of the phase space volume to be explored, in
529 order to allow for a finer scan of cost.

530 2.2.5. The Genetic Algorithm

531 A chromosome will be rendered as a concatenation of the options picked for every parameter, in the
532 order listed in the **.rng** file. They are generated on-the-fly, during the initialization phase of slave jobs,
533 and not in a centralized manner. When the pilot submits a slave job, it does not impose the parameter
534 configuration to be assessed by the slave, but expects the slave to generate such a configuration, based
535 on the past history of explored configurations, and assess its goodness. In other words, this GA is
536 *asynchronous*. Each so-far completed slave job will report the chromosome it has assessed, associated to
537 its fitness score, by concatenation to the end of the **done_so_far** file in the working directory. The pilot
538 script, periodically waking up from sleep to check the machine work load, will also verify whether new
539 entries were meanwhile added to **done_so_far**. If so, it will update the population of so-far best, non
540 redundant results in the working directory file **best_pop**, by sorting **done_so_far** by its fitness score,
541 then discarding all chromosomes that are very similar to slightly fitter essays and cutting the list off at
542 fitness levels of 70% of the best-so-far encountered fitness score (this "Darwinian selection" is the job of
543 awk script `$GACONF/chromdiv.awk`).

544 A new slave job will propose a novel parameter configuration by generating offspring of these "elite"
545 chromosomes in **best_pop**. In order to avoid reassessing an already seen, but not necessarily very
546 fit, configuration, it will also read the entire **done_so_far** file, in order to verify that the intended
547 configuration is not already in there. If so, genetic operators are again invoked, until a novel configuration
548 emerges. However, neither **best_pop** nor **done_so_far** do not yet include configurations that are
549 currently under evaluation on remote nodes. The asynchronous character of the GA does not allow
550 absolute guarantees that it will be a perfectly self-avoiding search, albeit measures have been taken in
551 this sense.

552 Upon lecture of **best_pop** and **done_so_far**, the parameter selector of the slave job (awk script
553 `$GACONF/make_childrenX.awk`) may randomly decide (with predefined probability) to perform
554 a "cross-over" between two randomly picked partners from **best_pop**, a "mutation" of a randomly
555 picked single parent, or a "spontaneous generation" (random initialization, ignoring the "ancestors" in
556 **best_pop**). Of course, if **best_pop** does not contain at least two individuals, the cross-over option is
557 not available, and if **best_pop** is empty, at the beginning of the simulation, then the mutation option is
558 unavailable as well: initially, all parameter configurations will be issued by "spontaneous generation".

559 2.2.6. Slave Processes – Startup, Learning Steps, Output Files and their Interpretation

560 Slave jobs run in temporary directories on the local file system of the nodes. A temporary directory
561 is being assigned an unambiguous attempt ID, appended to its default name "attempt". If the slave

562 job successfully completes, this attempt folder will be, after cleaning of temporary files, copied back
563 to the working directory. `$GACONF` also contains a sample working directory **D1-workdir**, displaying
564 typical results obtained from regression-based learning from **D1-datadir** – you may browse through it
565 in order to get acquainted with output files. The attempt ID associated to every chromosome is also
566 reported in the result file **done_so_far**. Inspect `$GACONF/D1-workdir/done_so_far`. Every
567 line thereof is formally divided by an equal sign in two sub-domains. The former encodes the actual
568 chromosome (field #1 stands for chosen DS, #2 for the descriptor scaling choice, ...). Fields at right
569 of the "=" report output data: attempt ID, fitting score and, as last field on line, the XV fitness score,
570 $\langle Q^2 \rangle - 2\sigma(Q^2)$ for regression, $\langle BA_{XV} \rangle - 2\sigma(BA_{XV})$ for classification problems as defined in
571 Introduction. The fitting score is calculated, for completeness, as the "mean-minus-two-sigma" of fitted
572 correlation coefficients and fitted balanced accuracy. Please do not allow yourself to be confused by
573 "fitting" (referring to statistics of the `svm-train` model fitting process, and concerning the instances
574 used to build the model), and Darwinian "fitness" (defined on the basis of XV results). The fitting score
575 is never used in the Darwinian selection process – it merely serves to inform the user about the loss of
576 accuracy between fitted and predicted/cross-validated property values.

577 **IMPORTANT!** Since every slave job will try to append its result line to **done_so_far** as soon as
578 it has completed calculations, chance may have different jobs on different nodes – each "seeing" the
579 working directory on a NFS-mounted partition – attempt to simultaneously write to **done_so_far**. This
580 may occasionally lead to corrupted lines, not matching the description above. Ignore them.

581 Yet, before results are copied back to the working directory, the actual work must be performed.
582 This begins by generating a chromosome, unless the `use_chromo` option instructs the job to apply an
583 externally imposed setup. After the chromosome is generated and stored in the current attempt folder,
584 it first must be interpreted. On one hand, the chromosome contains descriptor selection information.
585 The property-descriptor matrix matching the name reported in the first field of the chromosome, and
586 more precisely its scaled or non-scaled version, as indicated by the second field, will be the one
587 copied to the attempt folder, for further processing. The remaining chromosome elements must be
588 translated into the actual *libsvm* command line options. In this sense, the actual ϵ value for regression
589 is calculated by multiplying the `epsilon` parameter in the chromosome by the standard deviation of
590 the training property values. In this way, ϵ , representing 10 to 100% of the natural deviation of the
591 property, implicitly has the proper order of magnitude and proper units. Likewise, the chromosome
592 `gamma` parameter will be converted to the actual γ value, dividing by the mean vector dot product (for
593 kernel choices k1 or k3), or by the mean Euclidean distance (kernel choice k2), respectively. Since
594 the cost parameter is tricky to estimate, even in terms of magnitude orders, the chromosome features
595 a log-scale `cost` parameter, to be converted into the actual cost by taking the exponential thereof.
596 This, and also the stripping off of the "k" prefix in the kernel choice parameter, are also performed
597 at the chromosome interpretation stage, mode-dependently carried out by dedicated "decoder" `awk`
598 scripts **chromo2SVMreg.awk** and **chromo2SVMclass.awk**, respectively. At the end, the tool creates,
599 within the attempt folder, the explicit option line required to pilot `svm-train` according to the content
600 of the chromosome. This option line is stored in a one-line file **svm.pars** – check out, for example,
601 `$GACONF/D1-workdir/attempt.11118/svm.pars`, the option set that produced the fittest
602 models so-far.

603 A singled-out property-descriptor matrix and a command-line option set for `svm-train` are the
604 necessary and sufficient prerequisites to start XV. For each of the `ntrials` requested XV attempts, the
605 order of lines in the property-descriptor matrix is randomly changes. The reordered matrix is then split
606 into the requested `lo` folds. XV is explicit, and is not using the own facility of `svm-train`. Iteratively,
607 `svm-train` is used to learn a local model on the basis of all but the left-out fold. If adding of decoys
608 is desired, random decoy descriptor lines (of the same type, and having undergone the same scaling
609 strategy as training set descriptors) are added to make up 50% of the local learning set size (the `lo-1`
610 folds in use). These will differ at each successive XV attempt, if the pool of decoys out of which they
611 are drawn is much larger than the required half of training set size).

612 These local models are stored, and used to predict the entire training set, in which, however, the
613 identities of local "training" and locally left-out ("test") instances are known. These predictions are
614 separately reported into "train" and "test" files. The former report, for each instance, fitted property
615 values by local models having used the instance for training. The latter report predicted property values
616 by local models not having used it for training.

617 Quality of fit is checked on-the-fly, and failure to exceed a predefined fitting quality criterion triggers
618 a premature exit of the slave job, with a fictitious low fitness score associated to the current chromosome
619 – see the `fit_no_go` option. In such a case, the attempt subdirectory is deleted, and a new slave job is
620 launched instead to the defunct one.

621 Results in the attempt subdirectories of **D1-workdir** correspond to the default 12 times repeated
622 leave-1/3-out XV schemes. This means that $12 \times 3 = 36$ local models are generated. For each D1
623 ligand, there are exactly 12 of these local models that were not aware of that structure when they were
624 fitted, and 24 others that did benefit from its structure-property information upon building.

625 Files **final_train.pred.gz** in attempt subdirectories (provided as compressed .gz) report, for each
626 instance, the 24 "fitted" affinity values returned by models having used it for fitting (it is a 26-column
627 file, in which column #1 reports current numbering and #2 contains the experimental affinity
628 value). Reciprocally, the 14-column files **final_test.pred.gz** report the 12 prediction results stemming
629 from models not encountering instances at fitting stage. Furthermore, **consens_train.pred.gz** and
630 **consens_test.pred.gz** are "condensed" versions of the previous, in which the multiple predictions per
631 line have been condensed to their mean (column #2) and standard deviations (column #3), column #1
632 being the experimental property.

633 Eventually, the **stat** file found in each attempt subdirectory provides detailed statistics about fitting and
634 prediction errors in terms of root-mean-squared error, maximal error and determination coefficients. In a
635 classification process, correctly classified fractions and balanced accuracies are reported. Lines labeled
636 "local_train" compare every fitted value column from **final_train**, to the experimental data, whilst
637 "local_test" represent local model-specific XV results. By extracting all the lines "local_test:r_squared",
638 actually representing the XV coefficients of every local model, and calculating the mean and standard
639 deviations of these values, one may recalculate the fitness score of this attempt. The **stat** file lines
640 labeled "test" and "train" compare the consensus means as reported in the **consens_*.pred** files to
641 the experiment. Note: the determination coefficient "r_squared" between the mean of predictions
642 and experiment tends to exceed the mean of local determination coefficients, reporting individual

643 performances of local models. This is the "consensus" effect, not exploited in fitness score calculations
644 that rely on the latter mean of individual coefficients, penalized by twice their standard deviations.

645 Last but not least, this model building exercise had included an external prediction set of D5 ligands,
646 for which the current models were required to return a prediction of their D1 affinities. External
647 prediction is enabled as soon as model fitness exceeds the user-defined threshold `minlevel`. This
648 was the case for the best model corresponding to `$GACONF/D1-workdir/attempt.11118/`.
649 External prediction file names are a concatenation of external set name ("D5"), the locally used DS
650 ("treePH03_pruned"), the descriptor scaling status ("scaled") and the extension ".pred". As all these
651 external items are, by definition, "external" to all the $12 \times 3 = 36$ local models (no checking is performed
652 in order to detect potential overlaps with the training set), the 38-column external prediction file will
653 report the current numbering (column #1), the data reported in the first field of the .psvm files in column
654 #2 (here, the experimental D5 pK_i values), followed by 36 columns of individual predictions of the
655 modeled property (the D1 affinity values), by each of the local models. Since the external prediction set
656 .psvm files had been endowed with numeric data in the first field of each line, the present tool assumes
657 these to be corresponding experimental property values, to be compared to the predictions. Therefore, it
658 will take the consensus of the 36 individual D1 affinity predictions for each item, and compare them
659 to the experimental field. The results of external prediction challenge statistics are reported in the
660 **extval** file of the attempt subdirectory. First, the strictest comparison consists in calculating the RMSE
661 between experimental and predicted data, and the determination coefficient – these results are labeled
662 "Det" in the **extval** report. However, external prediction exercises may be challenging, and sometimes
663 useful even if the direct predicted-experimental error is large. For example, predicted values may not
664 quantitatively match experiment, but happen to be all offset by a common value or, in the weakest case,
665 nevertheless follow some linear trend – albeit of arbitrary slope and intercept. In order to test either of
666 these hypotheses, **extval** also reports statistics for (a) the predicted *vs.* experimental regression line at
667 fixed slope of 1.0, but with free intercept, labeled "FreeInt", and (b) the unconstrained, optimal linear
668 predicted-experimental correlation that may be established, labeled "Corr".

669 In this peculiar case, **extval** results seem unlikely to make any sense, because the tool mistakenly
670 takes the property data associated to external compounds (D5 affinities) for experimental D1 affinities,
671 to be confronted with their predicted alter-egos. Surprise – even the strict Determination statistics are
672 robustly positive on the fact that predicted D1 pK_i values quantitatively match experimental D5 pK_i
673 values. This is due to the very close biological relatedness of the two targets, which do happen to share
674 a lot of ligands. Indeed, 25% of the instances of the D5 external set were also reported ligands of D1,
675 and, as such, part of the model training set. Yet, 3 out of 4 D5 ligands were genuinely new. This
676 notwithstanding, their predicted D1 affinities came close to observed D5 values – a quite meaningful
677 result confirming the extrapolative prediction abilities of the herein built model.

678 At this point, local SVM models and other temporary files are deleted, the chromosome associated to
679 it attempt ID, fitting and fitness scores is being appended to **done_so_far**, and the attempt subfolder now
680 containing only setup information and results is moved from its temporary location on the cluster back
681 to the working directory (with exception of the workstation-based implementation, when temporary and
682 final attempt subdirectory location are identical). The slave job successfully exits.

683 2.2.7. Reconstruction and Off-package Use of Optimally Parameterized *libsvm* Models

684 As soon as the number of processed attempts reaches `maxconfigs`, the pilot script stops
 685 rescheduling new slave jobs. Ongoing slave processes are allowed to complete – therefore, the final
 686 number of lines in **done_so_far** may slightly exceed `maxconfigs`. Before exiting, the pilot will also
 687 clean the working directories, by deleting all the attempt sub-folders corresponding to the less successful
 688 parameterization schemes, which did not pass the selection hurdle and are not listed in **best_pop**. A
 689 trace of their statistical parameters is saved in a directory called "losers".

690 The winning attempt sub-folders contain all the fitted and cross-validated/predicted property tables
 691 and statistics reports, but not the battery of local SVM models that served to generate them. These
 692 were not kept, because they may be quite bulky files which may not be allowed to accumulate in large
 693 numbers on the disk (at the end of the run, before being able to decide which attempts can be safely
 694 discarded, there should have been `maxconfigs × lo times ntrials` of them). However, they – or,
 695 actually, equivalent (recall that local model files are tributary to the random regrouping of kept *vs.* left-out
 696 instances at each XV step) – model files can be rebuilt (and kept) by rerunning the pilot script with
 697 the `use_chromo` option pointing to a file with chromosomes encoding the desired parameterization
 698 scheme(s). If `use_chromo` points to a multi-line file, the rebuilding process may be parallelized: as the
 699 slave jobs proceed, new attempt sub-folders – now each containing `lo times ntrials libsvm` model
 700 files – will appear in the working directory (but not necessarily in the listing order of `use_chromo`).

701 **IMPORTANT!** The file of preferred setup chromosomes should only contain the parameter fields,
 702 but not the chromosome evaluation information. Therefore, you cannot use **best_pop** *per se*, or ``head`
 703 `-top`` of **best_pop** as `use_chromo` file, unless you first remove the right-most fields, starting at the
 704 "equal" separator:

```
705
706 head -top best_pop | sed 's/=.*//' > my_favorite_chromosomes.lst
707 $GACONF/pilot_<scheme>.csh workdir=<working directory used for GA
708 run> use_chromo=my_favorite_chromosomes.lst
```

709
 710 Randomness at learning/left-out splitting stages may likely cause the fitness score upon refitting to
 711 (slightly) differ from the initial one (on the basis of which the setup was considered interesting). If the
 712 difference is significant, it means that the XV strategy was not robust enough – *i.e.* was not repeated
 713 sufficiently often (increase `ntrials`).

714 Local model files can now be used for consensus predictions, and the degree of divergence of
 715 their prediction for an external instance may serve as prediction confidence score [5]. Remember that
 716 `svm-predict` using those model files should be called on preprocessed descriptor files, as outlined in
 717 the Command-Line Launch subsection.

718 Alternatively, you may rebuild your *libsvm* models manually, using the `svm-train` parameter
 719 command line written out in successful attempt sub-folders. In this case, you may tentatively rebuild
 720 the models on your brute descriptor files, or on `svm-scale-preprocessed` descriptor files – unless the
 721 successful recipe requires pruning.

722 3. Conclusions

723 This parameter configurator for *libsvm* addresses several important aspects in SVM model building,
724 in view – but not exclusively dedicated to – chemoinformatics applications.

725 First, the approach co-opts an important decision making of the model building process: the choice
726 of descriptors/attributes, and their best preprocessing strategies, into the optimization procedure. This is
727 important, because the employed descriptor space is a primordial determinant of modeling success, and
728 determines the optimal operational parameters of *libsvm*.

729 Next, an aggressive, repeated cross-validation scheme, introducing a penalty proportional to the
730 fluctuation of cross-validation propensity on the local grouping of learning *vs.* left-out instances is used
731 to select an operational parameter set leading to a model of maximal robustness, not to a model owing its
732 apparent quality to a lucky ordering of instances in the training set. This fitness score is in our opinion a
733 good estimator of the extrapolative predictive power of the model.

734 Furthermore, the approach allows decoy instances to be added in order to expand learned problem
735 space zone, without however adding uncertainty to the reported statistics (statistics of decoy-free and
736 decoy-based models are directly comparable, because they both focus on the actual training instances
737 and their confirmed experimental properties, ignoring "presumed" inactive decoy property values).

738 Last but not least, the approach is versatile, covering both regression and classification problems and
739 supporting a large variety of parallel deployment schemes. It was, for example, successfully used to
740 generate very large (> 9000 instances) dataset-based chemogenomics models [6].

741 **BA – Balanced Accuracy (of classification, the mean of sensitivity and specificity) DS – Descriptor**
742 **Space, GA – Genetic Algorithm, SVM – Support Vector Machine, SVR – Support Vector**
743 **Regression, SVC – Support Vector Classification, RMSE – Root Mean Squared Error, XV -**
744 **Cross-validation.**

745 **The authors thank the High Performance Computing centers of the Universities of Strasbourg,**
746 **France, and Cluj, Romania, for having hosted a part of the herein reported development work. This**
747 **work was also supported in part by a Grant-in-Aid for Young Scientists from the Japanese Society**
748 **for the Promotion of Science (Kakenhi (B) 25870336). This research was additionally supported**
749 **by the Funding Program for Next Generation World-Leading Researchers as well as the CREST**
750 **program of the Japan Science and Technology Agency.**

751 References

- 752 1. Chang, C.C.; Lin, C.J. LIBSVM: A library for support vector machines. *ACM Transactions on*
753 *Intelligent Systems and Technology* **2011**, *2*, 27:1–27.
- 754 2. Gaulton, A.; Bellis, L.J.; Bento, A.P.; Chambers, J.; Davies, M.; Hersey, A.; Light, Y.;
755 McGlinchey, S.; Michalovich, D.; Al-Lazikani, B.; Overington, J.P. ChEMBL: a large-scale
756 bioactivity database for drug discovery. *Nucl. Ac. Res.* **2011**, *40*, D1100–D1107.

- 757 3. Ruggiu, F.; Marcou, G.; Varnek, A.; Horvath, D. Isida Property-labelled Fragment Descriptors.
758 *Molecular Informatics* **2010**, *29*, 855–868.
- 759 4. Bonachera, F.; Parent, B.; Barbosa, F.; Froloff, N.; Horvath, D. Fuzzy Tricentric Pharmacophore
760 Fingerprints. 1 - Topological Fuzzy Pharmacophore Triplets and adapted Molecular Similarity
761 Scoring Schemes. *J. Chem. Inf. Model.* **2006**, *46*, 2457–2477.
- 762 5. Horvath, D.; Marcou, G.; Varnek, A. Predicting the Predictability: A Unified Approach to the
763 Applicability Domain Problem of QSAR Models. *J. Chem. Inf. Model.* **2009**, *49*, 1762–1776.
- 764 6. Brown, J.B.; Okuno, Y.; Marcou, G.; Varnek, A.; Horvath, D. Computational chemogenomics:
765 Is it more than inductive transfer? *J. Comput.-Aided Mol. Des.* **2014**, *28*, 597–618.

766 © August 7, 2014 by the authors; submitted to *Challenges* for possible open access
767 publication under the terms and conditions of the Creative Commons Attribution license
768 <http://creativecommons.org/licenses/by/3.0/>.