

ParamILS: Iterated Local Search in Parameter Configuration Space

Jiří Kubalík
Department of Cybernetics, CTU Prague

Substantial part of this material is based on the paper
Frank Hutter, Holger H. Hoos, Kevin Leyton-Brown, and Thomas Stützle: ParamILS: An Automatic Algorithm Configuration Framework,
Journal of Artificial Intelligence Research (JAIR),
volume 36, pp. 267-306, October 2009.
See <http://www.cs.ubc.ca/~hutter/papers/Hutter09PhD.pdf>



<http://cw.felk.cvut.cz/doku.php/courses/a0m33eoa/start>

Algorithm Configuration (Parameter Tuning) Problem

The objective of the parameter configuration (parameter tuning) problem is to find the parameter configuration $\theta \in \Theta$ resulting in the best performance of \mathcal{A} on distribution \mathcal{D} .

There are many ways of measuring the **cost**, $o(\mathcal{A}, \theta, I, s)$, of running algorithm \mathcal{A} with parameter configuration θ on an instance I , using seed s in case of randomized algorithm. For example, we might be interested in

- the computational resources consumed by the given algorithm (such as runtime, memory or communication bandwidth), or
- the approximation error, or
- the improvement achieved over an instance-specific reference cost,
- the quality of the solution found.

Algorithm Configuration (Parameter Tuning) Problem

The behaviour of the algorithms can vary significantly between multiple runs on different instances or when randomized algorithms are run repeatedly with fixed parameters on a single problem instance.

Therefore, the cost of a candidate solution θ is defined as

$$c(\theta) = m(O_\theta)$$

the statistical population parameter m of the cost distribution $O_\theta(\mathcal{A}, \theta, \mathcal{D})$, over instances drawn from distribution of instances, \mathcal{D} , and multiple independent runs.

An optimal solution, θ^* , minimizes $c(\theta)$:

$$\theta^* \in \arg \min_{\theta \in \Theta} c(\theta).$$

For example, we might aim to minimize **mean runtime** or **median solution cost**.

The $O_\theta(\mathcal{A}, \theta, \mathcal{D})$ is typically unknown, so we can only acquire approximations of their statistics, $c(\theta)$, based on a limited number of samples (i.e. the cost of single executions of $\mathcal{A}(\theta)$) – let's denote an approximation of $c(\theta)$ based on N samples by $\hat{c}_N(\theta)$.

- For deterministic algorithms, the algorithm \mathcal{A} is run on $N \leq M$ instances (M is the size of the finite training set of instances).
- For randomized, algorithms, we can run multiple runs with different seeds if $M < N$.

Searching the Parameter Configuration Space

Manually-executed local search in parameter configuration space

1. Begin with some initial parameter configuration.
2. Experiment with modifications to single parameter values at a time, accepting new configurations whenever they result in improved performance (**iterative first improvement procedure**).
3. Repeat step 2 until no single-parameter change yields an improvement.

This procedure stops as soon as it **reaches a local optimum** (a parameter configuration that cannot be improved by modifying a single parameter value).

Issues when trying to automate the process

- Which parameter configurations should be evaluated?
- Which problem instances should be used and how many runs should be performed on each instance?
Stochastic nature of the algorithm configuration problem.
- Which cutoff time κ_i (the maximum amount of time the configured algorithm is allowed to use) should be used for each run?

ParamILS: Iterated Local Search in Parameter Configuration Space

Employs **Iterated Local Search** that builds a chain of local optima by iterating through a main loop consisting of:

1. a solution perturbation to escape from local optima,
2. a subsidiary local search procedure and
3. an acceptance criterion to decide whether to keep or reject a newly obtained candidate solution.

$ParamILS(\theta_0, r, p_{restart}, s)$

1. uses a combination of default and random settings for initialization,
 θ_0 is the initial parameter configuration, and
 r is the number of randomly chosen configurations for initialization,
2. uses a one-exchange neighborhood (one parameter is modified in each search step),
3. employs iterative first improvement as a subsidiary local search procedure,
4. uses a fixed number, s , of random moves for perturbation,
5. always accepts better or equally-good parameter configurations,
6. re-initializes the search at random with probability $p_{restart}$.

ParamILS(N): Algorithm

```

Input : Initial configuration  $\theta_0 \in \Theta$ , algorithm parameters  $r, p_{restart}$ , and  $s$ .
Output : Best parameter configuration  $\theta$  found.

1 for  $i = 1, \dots, r$  do
2    $\theta \leftarrow \text{random } \theta \in \Theta$ ;
3   if  $\text{better}(\theta, \theta_0)$  then  $\theta_0 \leftarrow \theta$ ;
4  $\theta_{ils} \leftarrow \text{IterativeFirstImprovement}(\theta_0)$ ; First iteration of the ILS procedure
5 while not  $\text{TerminationCriterion}()$  do
6    $\theta \leftarrow \theta_{ils}$ ;
7   // ===== Perturbation
8   for  $i = 1, \dots, s$  do  $\theta \leftarrow \text{random } \theta' \in \text{Nbh}(\theta)$ ;
9   // ===== Basic local search
10   $\theta \leftarrow \text{IterativeFirstImprovement}(\theta)$ ;
11  // ===== AcceptanceCriterion
12  if  $\text{better}(\theta, \theta_{ils})$  then  $\theta_{ils} \leftarrow \theta$ ;
13  with probability  $p_{restart}$  do  $\theta_{ils} \leftarrow \text{random } \theta \in \Theta$ ;
14 return overall best  $\theta_{inc}$  found;
15 Procedure  $\text{IterativeFirstImprovement}(\theta)$ 
16 repeat
17    $\theta' \leftarrow \theta$ ;
18   foreach  $\theta'' \in \text{Nbh}(\theta')$  in randomized order do
19     if  $\text{better}(\theta'', \theta')$  then  $\theta \leftarrow \theta''$ ; break;
20 until  $\theta' = \theta$ ;
21 return  $\theta$ ;

```

Initialization

Main body of the ILS procedure

How to define $\text{better}(\theta'', \theta')$?

BasicILS(N): Procedure $better_N(\theta_1, \theta_2)$

Basic variant, *BasicILS*(N), uses procedure $better_N(\theta_1, \theta_2)$ that compares two cost approximations $\hat{c}_N(\theta_1)$ and $\hat{c}_N(\theta_2)$ based on exactly N samples from the respective cost distributions $O_\theta(\mathcal{A}, \theta_1, \mathcal{D})$ and $O_\theta(\mathcal{A}, \theta_2, \mathcal{D})$ – the same N instances are used for all configurations θ_i .

Procedure $better_N(\theta_1, \theta_2)$ simply compares estimates $\hat{c}_N(\theta_1)$ and $\hat{c}_N(\theta_2)$ based on the same N instances using the same random seeds.

- It updates the best-so-far solution, θ_{inc} .

Input	: Parameter configuration θ_1 , parameter configuration θ_2
Output	: True if θ_1 does better than or equal to θ_2 on the first N instances; false otherwise
Side Effect	: Adds runs to the global caches of performed algorithm runs \mathbf{R}_{θ_1} and \mathbf{R}_{θ_2} ; potentially updates the incumbent θ_{inc}
1	$\hat{c}_N(\theta_2) \leftarrow objective(\theta_2, N)$
2	$\hat{c}_N(\theta_1) \leftarrow objective(\theta_1, N)$
3	return $\hat{c}_N(\theta_1) \leq \hat{c}_N(\theta_2)$

FocusedILS

The question is how to choose the optimal number of training instances, N ?

- Using too small N leads to good training performance, but poor generalization to previously unseen test benchmarks.
- On the other hand, we cannot evaluate every parameter configuration on an enormous training set - if we did, search progress would be unreasonably slow.

FocusedILS is a variant of ParamILS that **adaptively varies the number of training samples** considered from one parameter configuration to another in order **to focus samples on promising configurations**.

- $N(\theta)$ denotes the number of runs available to estimate the cost statistic $c(\theta)$.

The question is how to compare two parameter configurations θ_1 and θ_2 for which $N(\theta_1) \leq N(\theta_2)$?

- *What if we computed the empirical statistics based on the available number of runs for each configuration?*

FocusedILS

The question is how to choose the optimal number of training instances, N ?

- Using too small N leads to good training performance, but poor generalization to previously unseen test benchmarks.
- On the other hand, we cannot evaluate every parameter configuration on an enormous training set - if we did, search progress would be unreasonably slow.

FocusedILS is a variant of ParamILS that **adaptively varies the number of training samples** considered from one parameter configuration to another in order **to focus samples on promising configurations**.

- $N(\theta)$ denotes the number of runs available to estimate the cost statistic $c(\theta)$.

The question is how to compare two parameter configurations θ_1 and θ_2 for which $N(\theta_1) \leq N(\theta_2)$?

- *What if we computed the empirical statistics based on the available number of runs for each configuration?*

Can lead to systematic bias if, for example, the first instances are easier than the average ones.

FocusedILS: Domination

Domination: Configuration θ_1 dominates θ_2 when at least as many runs have been conducted on θ_1 as on θ_2 , and the performance of $\mathcal{A}(\theta)$ on the first $N(\theta_1)$ runs is at least as good as that of $\mathcal{A}(\theta_2)$ on all of its runs.

θ_1 dominates θ_2 if and only if $N(\theta_1) \geq N(\theta_2)$ and $\hat{c}_{N(\theta_2)}(\theta_1) \leq \hat{c}_{N(\theta_2)}(\theta_2)$.

ParamILS version, called FocusedILS, encodes a comparison strategy based on the domination in procedure $\textit{better}_{\textit{Foc}}(\theta_1, \theta_2)$.

FocusedILS: Procedure $better_{Foc}(\theta_1, \theta_2)$

```

Input      : Parameter configuration  $\theta_1$ , parameter configuration  $\theta_2$ 
Output    : True if  $\theta_1$  dominates  $\theta_2$ , false otherwise
Side Effect: Adds runs to the global caches of performed algorithm runs  $\mathbf{R}_{\theta_1}$  and  $\mathbf{R}_{\theta_2}$ ; updates the
               global counter  $B$  of bonus runs, and potentially the incumbent  $\theta_{inc}$ 
1  $B \leftarrow B + 1$ 
2 if  $N(\theta_1) \leq N(\theta_2)$  then
3    $\theta_{min} \leftarrow \theta_1; \theta_{max} \leftarrow \theta_2$ 
4   if  $N(\theta_1) = N(\theta_2)$  then  $B \leftarrow B + 1$ 
5 else  $\theta_{min} \leftarrow \theta_2; \theta_{max} \leftarrow \theta_1$ 
6 repeat
7    $i \leftarrow N(\theta_{min}) + 1$ 
8    $\hat{c}_i(\theta_{max}) \leftarrow objective(\theta_{max}, i)$  // If  $N(\theta_{min}) = N(\theta_{max})$ , adds a new run to  $\mathbf{R}_{\theta_{max}}$ .
9    $\hat{c}_i(\theta_{min}) \leftarrow objective(\theta_{min}, i)$  // Adds a new run to  $\mathbf{R}_{\theta_{min}}$ .
10 until  $dominates(\theta_1, \theta_2)$  or  $dominates(\theta_2, \theta_1)$ 
11 if  $dominates(\theta_1, \theta_2)$  then
12   // ===== Perform  $B$  bonus runs.
13    $\hat{c}_{N(\theta_1)+B}(\theta_1) \leftarrow objective(\theta_1, N(\theta_1) + B)$  // Adds  $B$  new runs to  $\mathbf{R}_{\theta_1}$ .
14    $B \leftarrow 0$ 
15   return true
16 else return false
17 Procedure  $dominates(\theta_1, \theta_2)$ 
18 if  $N(\theta_1) < N(\theta_2)$  then return false
19 return  $objective(\theta_1, N(\theta_2)) \leq objective(\theta_2, N(\theta_2))$ 

```


Adaptive Capping of Algorithm Runs

Often, the computational resources are wasted with evaluating a parameter configuration that is much worse than other, previously-seen configurations.

Ex.: Let's assume a case where parameter configuration θ_1 takes a total of 10 seconds to solve $N = 100$ instances (i.e. it has a mean runtime of 0.1 seconds per instance), and another parameter configuration θ_2 takes 100 seconds to solve the first of these instances.

Clearly, when comparing the mean runtimes of θ_1 and θ_2 based on this set of instances, it is not necessary to run θ_2 on remaining 99 instances. Instead, we can terminate the first run of θ_2 after $10 + \epsilon$ seconds, which is a lower bound on θ_2 's mean runtime of $0.1 + \epsilon/100$. This lower bound exceeds the mean runtime of θ_1 , so we can already be sure that θ_2 cannot do better than θ_1 .

Question is how to determine the cutoff time for each run of the target algorithm, \mathcal{A} , in an automated way?

Adaptive capping is based on the idea of avoiding unnecessary runs of the algorithm \mathcal{A} by **developing bounds on the performance measure to be optimized.**

- **Trajectory-preserving capping** – provably does not change *BasicILS*'s search trajectory, but can lead to large computational savings.
- **Aggressive capping** – potentially yielding even better performance.

Trajectory-Preserving Capping: Summary

The *trajectory-preserving capping* computes the upper bound on the cumulative runtime **from the best configuration encountered in the current ILS iteration**.

New ILS iteration starts with a **perturbation of the best-so-far configuration θ** .

Frequently, the new parameter configuration θ is of poor quality, thus **the capping criterion does not apply as quickly as it could** if the comparison was performed against the overall *incumbent*.

FocusedILS: Adaptive Capping

FocusedILS varies the number of runs used to evaluate each parameter configuration.

How can the adaptive capping be used here?

Recommended Material

Frank Hutter, Holger H. Hoos, Kevin Leyton-Brown, and Thomas Stützle: ParamILS: An Automatic Algorithm Configuration Framework. In *Journal of Artificial Intelligence Research (JAIR)*, volume 36, pp. 267-306, October 2009.

Other papers and SW available at <http://www.cs.ubc.ca/labs/beta/Projects/ParamILS/>

