

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

Let \mathcal{A} denote an algorithm, whose parameters are to be optimized for a **probability distribution of problem instances**, \mathcal{D} .

\mathcal{D} may be given

- implicitly, as through a random instance generator or a distribution over such generators, or
- as the uniform distribution over a finite sample of problem instances.

With each of the algorithm parameters, $p_1 \dots p_k$, a domain of possible values is associated and the parameter configuration space, Θ , is the cross-product of these domains (or a subset thereof).

- We assume that all parameter domains are finite sets.
This assumption can be met by discretizing all numerical parameters to a finite number of values.
- While parameters may be ordered, we do not exploit such ordering relations \implies all parameters are finite and categorical.

The elements of Θ are called parameter configurations, θ_i , and $\mathcal{A}(\theta)$ denotes algorithm \mathcal{A} with parameter configuration $\theta \in \Theta$.

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 $M < N$.

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 $p_{restart}$.

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} .

<p>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}</p>

- 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 $c(\theta)$.

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

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 $\hat{c}_{N(\theta_2)}(\theta_1) \leq \hat{c}_{N(\theta_2)}(\theta_2)$.

FocusedILS – procedure $better_{Foc}(\theta_1, \theta_2)$ implements a comparison strategy based on the domination

1. first it acquires one additional run for the configuration i having smaller $N(\theta_i)$, or one run for both configurations if $N(\theta_1) = N(\theta_2)$;
2. then, it continues performing runs in this way until one configuration dominates the other.
It returns true if θ_1 dominates θ_2 , and false otherwise.

It keeps track of the total number, B , of configurations evaluated since the last improving step.

- Whenever $better_{Foc}(\theta_1, \theta_2)$ returns true, B extra (bonus) runs are performed for θ_1 and B is reset to 0.
- This way it is ensured that many runs are performed with good configurations \implies the error made in every comparison of two configurations θ_1 and θ_2 decreases on expectation.

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/>

