

# Information Utilization Ratio in Heuristic Optimization Algorithms

Junzhi Li and Ying Tan  
Computational Intelligence Laboratory  
Peking University  
{ljz,ytan}@pku.edu.cn

September 11, 2018

## Abstract

Heuristic algorithms are able to optimize objective functions efficiently because they use intelligently the information about the objective functions. Thus, information utilization is critical to the performance of heuristics. However, the concept of information utilization has remained vague and abstract because there is no reliable metric to reflect the extent to which the information about the objective function is utilized by heuristic algorithms. In this paper, the metric of information utilization ratio (IUR) is defined, which is the ratio of the utilized information quantity over the acquired information quantity in the search process. The IUR proves to be well-defined. Several examples of typical heuristic algorithms are given to demonstrate the procedure of calculating the IUR. Empirical evidences on the correlation between the IUR and the performance of a heuristic are also provided. The IUR can be an index of how finely an algorithm is designed and guide the invention of new heuristics and the improvement of existing ones.

## 1 Introduction

In the field of computer science, many heuristic algorithms have been developed to solve complex non-convex optimization problems. Although optimal solutions are not guaranteed to be found, heuristics can often find acceptable solutions at affordable cost. The key to designing a heuristic algorithm is to use heuristic information about the objective function. Many algorithms [1, 2, 3] are claimed to be reasonably designed because they use heuristic information intelligently. Even more algorithmic improvement works [4, 5, 6] are claimed to be significant because they use more heuristic information or use heuristic information more thoroughly than the original algorithms.

Empirically, heuristic information is used more thoroughly in more advanced algorithms. Suppose there are two search algorithms A and B for one dimen-

sional optimization. Algorithm A compares the evaluation values of the solutions  $x_1$  and  $x_2$  to decide which direction (left or right) is more promising, while algorithm B uses their evaluation values to calculate both the direction and the step size for the next search. If the underlying distribution of objective functions is already known, then algorithm B is able to search faster than algorithm A if they are both reasonably designed because more information is utilized by algorithm B. It has been a common sense in the field of heuristic search that the extent of information utilization in a heuristic algorithm is crucial to its performance.

However, so far there is no reliable metric to reflect the extent of information utilization because unlike direct performance analyses [7, 8], this issue seems abstract. Especially, it is very difficult to measure how much information is used by an optimization algorithm.

In this paper, based on some basic concepts in the information theory, a formal definition of the information utilization ratio (IUR) is proposed, which is defined as the ratio of the utilized information quantity over the acquired information quantity in the search process. It is shown theoretically that IUR is well-defined. Examples of typical heuristic algorithms are also given to demonstrate the procedure of calculating IURs.

Theoretically, IUR itself is a useful index of how finely an algorithm is designed, but we still expect it to be practically serviceable, that is, we need to study the correlation between IUR and performance. However, the correlation between IUR and performance of heuristics is not so straightforward as some may expect. The performance of an optimization algorithm depends not only on the extent of information utilization but also on the manner of information utilization. Still, for algorithms that utilize information in similar manners, the influence of the IUR is often crucial, as is illustrated in the experiments.

After all, the metric of IUR helps researchers construct a clear (but not deterministic) relationship between the design and the performance of an optimization algorithm, which makes it possible that researchers can to some extent predict the performance of an algorithm even before running it. Thus, the IUR can be a useful index for guiding the design and the improvement of heuristic optimization algorithms.

## 2 Information Utilization Ratio

**Definition 1** (Information Entropy). *The information entropy of a discrete random variable  $X$  with possible values  $x_i$  and probability density  $p(x_i)$  is defined as follows.*

$$H(X) = - \sum_i p(x_i) \log p(x_i). \quad (1)$$

**Definition 2** (Conditional Entropy). *The conditional entropy of two discrete random variables  $X$  and  $Y$  with possible values  $x_i$  and  $y_j$  respectively and joint*

probability density  $p(x_i, y_j)$  is defined as follows.

$$H(X|Y) = - \sum_{i,j} p(x_i, y_j) \log \frac{p(x_i, y_j)}{p(y_j)}. \quad (2)$$

Some elementary properties of information entropy and conditional entropy are frequently used in this paper, which however cannot be present here due to the limitation of space. We refer readers who are unfamiliar with the information theory to the original paper [9] or other tutorials.

The following lemma defines a useful function for calculating the IURs of various algorithms.

**Lemma 1.** *If  $\eta_1, \eta_2, \dots, \eta_{g+1} \in \mathbb{R}$  are independent identically distributed random variables,*

$$\begin{aligned} & H(I(\min(\eta_1, \eta_2, \dots, \eta_g) < \eta_{g+1})) \\ &= -\frac{g}{g+1} \log \frac{g}{g+1} - \frac{1}{g+1} \log \frac{1}{g+1} \triangleq \pi(g), \end{aligned} \quad (3)$$

where  $I(x < y) = \begin{cases} 1 & \text{if } x < y \\ 0 & \text{otherwise} \end{cases}$  is the indicator function.

$\pi(g) \in (0, 1]$  is a monotonic decreasing function of  $g$ .

**Definition 3** (Objective Function). *The objective function is a mapping  $f : \mathcal{X} \mapsto \mathcal{Y}$ , where  $\mathcal{Y}$  is a totally ordered set.*

$\mathcal{X}$  is called the search space. The target of an optimization algorithm is to find a solution  $x \in \mathcal{X}$  with the best evaluation value  $f(x) \in \mathcal{Y}$ .

**Definition 4** (Optimization Algorithm). *An optimization algorithm  $\mathcal{A}$  is defined as follows.*

---

**Algorithm 1** Optimization Algorithm  $\mathcal{A}$

---

- 1:  $i \leftarrow 0$ .
  - 2:  $D_0 \leftarrow \emptyset$ .
  - 3: **repeat**
  - 4:    $i \leftarrow i + 1$ .
  - 5:   Sample  $X_i \in 2^{\mathcal{X}}$  with distribution  $\mathcal{A}_i(D_{i-1})$ .
  - 6:   Evaluate  $f(X_i) = \{f(x) | x \in X_i\}$ .
  - 7:    $D_i \leftarrow D_{i-1} \cup \bigcup_{x \in X_i} \{x, f(x)\}$ .
  - 8: **until**  $i = g$ .
- 

*In each iteration,  $\mathcal{A}_i$  is a mapping from  $2^{\mathcal{X} \times \mathcal{Y}}$  to the set of all distributions over  $2^{\mathcal{X}}$ .  $\mathcal{A}_1(D_0)$  is a pre-fixed distribution for sampling solutions in the first iteration.  $g$  is the maximal iteration number.*

In each iteration, the input of the algorithm  $D_{i-1}$  is the historical information, which is a subset of  $\mathcal{X} \times \mathcal{Y}$ , and the output  $\mathcal{A}_i(D_{i-1})$  is a distribution over  $2^{\mathcal{X}}$ , with which the solutions to be evaluated next are drawn. Note that the output  $\mathcal{A}_i(D_{i-1})$  is deterministic given  $D_{i-1}$ .

By randomizing the evaluation step (consider  $y = f(x)$  as a random variable), we are able to investigate how much acquired information is used in an optimization algorithm. That is, to what extent the action of the algorithm will change when the acquired information changes. Review the example in the introduction. It is clear that the algorithm A only uses the information of "which one is better", while the information of evaluation values are fully utilized by the algorithm B. But how to express such an observation? Any change in  $y_1$  or  $y_2$  would cause the algorithm B to search a different location, while only when  $I(y_1 > y_2)$  changes would the action of the algorithm A change. So, the quantity of utilized information can be expressed by the information entropy of an algorithm's action. The entropy of the action of the algorithm A is one bit, while the entropy of the action of the algorithm B is equal to the entropy of the evaluation values. Assume  $Z$  is the "action" of the algorithm,  $X$  is the positions of the solutions,  $Y$  is the evaluation values, (they are all random variables), then we can roughly think the information utilization ratio is  $H(Z|X)/H(Y|X)$ . However, optimization algorithms are iterative processes, so the formal definition is more complicated.

**Definition 5** (Information Utilization Ratio). *If  $\mathcal{A}$  is an optimization algorithm, the information utilization ratio of  $\mathcal{A}$  is defined as follows.*

$$IUR_{\mathcal{A}}(g) = \frac{\sum_{i=1}^g H(Z_i|\bar{X}_{i-1}, \bar{Z}_{i-1})}{\sum_{i=1}^g H(Y_i|\bar{X}_i, \bar{Y}_{i-1})}, \quad (4)$$

where  $g$  is the maximal iteration number,  $X = \{X_1, X_2, \dots, X_g\}$  is the set of all sets of evaluated solutions,  $Y = \{f(X_1), f(X_2), \dots, f(X_g)\}$  is the set of all sets of evaluation values,  $Z = \{\mathcal{A}_1(D_0), \mathcal{A}_2(D_1), \dots, \mathcal{A}_g(D_{g-1})\}$  is the output distributions in all iterations of algorithm  $\mathcal{A}$ ,  $\bar{X}_i \triangleq \{X_1, \dots, X_i\}$ ,  $\bar{Y}_i \triangleq \{Y_1, \dots, Y_i\}$ ,  $\bar{Z}_i \triangleq \{Z_1, \dots, Z_i\}$ ,  $\bar{X}_0 = \bar{Y}_0 = \bar{Z}_0 = \emptyset$ .

Fig. 1 shows the relationship among these random variables. Generally,  $X_i$  is acquired by sampling with the distribution  $Z_i$ ,  $Y_i$  is acquired by evaluating  $X_i$ , and  $Z_i$  is determined by the algorithm according to the historical information  $\bar{X}_{i-1}$  and  $\bar{Y}_{i-1}$ .

For deterministic algorithms (i.e.,  $H(X_i|Z_i) = 0$ ), the numerator degenerates to  $H(Z)$ . If function evaluations are independent, the denominator degenerates to  $\sum_{i=1}^g H(Y_i|X_i)$ .

The following theorem guarantees that IUR is well defined.

**Theorem 1.** *If  $0 < \sum_{i=1}^g H(Y_i|\bar{X}_i, \bar{Y}_{i-1}) < \infty$ , then  $0 \leq IUR_{\mathcal{A}}(g) \leq 1$ .*

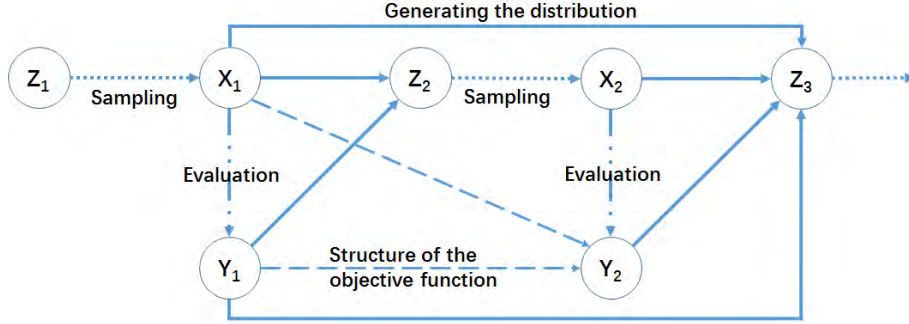


Figure 1: Graphic Model

*Proof.*

$$H(X, Z) - \sum_{i=1}^g H(X_i | Z_i) \quad (5)$$

$$= \sum_{i=1}^g H(X_i, Z_i | \bar{X}_{i-1}, \bar{Z}_{i-1}) - \sum_{i=1}^g H(X_i | \bar{X}_{i-1}, \bar{Z}_i) \quad (6)$$

$$= \sum_{i=1}^g H(Z_i | \bar{X}_{i-1}, \bar{Z}_{i-1}) \quad (7)$$

$$= \sum_{i=2}^g H(\bar{Z}_i | \bar{X}_{i-1}) - \sum_{i=2}^g H(\bar{Z}_{i-1} | \bar{X}_{i-1}) \quad (8)$$

$$= \sum_{i=2}^g H(\bar{Z}_i | \bar{X}_{i-1}) - \sum_{i=2}^g H(\bar{Z}_i | \bar{X}_{i-1}, \bar{Y}_{i-1}) - \sum_{i=2}^g H(\bar{Z}_{i-1} | \bar{X}_{i-1}) + \sum_{i=2}^g H(\bar{Z}_{i-1} | \bar{X}_{i-1}, \bar{Y}_{i-2}) \quad (9)$$

$$= \sum_{i=2}^g -H(\bar{Y}_{i-1} | \bar{Z}_i, \bar{X}_{i-1}) + \sum_{i=2}^g H(\bar{Y}_{i-1} | \bar{X}_{i-1}) + \sum_{i=2}^g H(\bar{Y}_{i-2} | \bar{Z}_{i-1}, \bar{X}_{i-1}) - \sum_{i=2}^g H(\bar{Y}_{i-2} | \bar{X}_{i-1}) \quad (10)$$

$$= \sum_{i=2}^g -H(\bar{Y}_{i-1} | \bar{Z}_i, \bar{X}_{i-1}) + \sum_{i=2}^g H(\bar{Y}_{i-2} | \bar{Z}_{i-1}, \bar{X}_{i-2}) + \sum_{i=2}^g H(Y_{i-1} | \bar{X}_{i-1}, \bar{Y}_{i-2}) \quad (11)$$

$$= -H(\bar{Y}_{g-1} | \bar{Z}_g, \bar{X}_{g-1}) + \sum_{i=1}^g H(Y_i | \bar{X}_i, \bar{Y}_{i-1}) - H(Y_g | \bar{X}_g, \bar{Y}_{g-1}) \quad (12)$$

$$\leq \sum_{i=1}^g H(Y_i | \bar{X}_i, \bar{Y}_{i-1}). \quad (13)$$

Eq. (8) holds because

$$H(Z_1) = 0. \quad (14)$$

Eq. (9) holds because

$$H(\bar{Z}_i|\bar{X}_{i-1}, \bar{Y}_{i-1}) = H(\bar{Z}_{i-1}|\bar{X}_{i-1}, \bar{Y}_{i-2}) = 0. \quad (15)$$

Eq. (11) holds because

$$\sum_{i=2}^g H(\bar{Y}_{i-2}|\bar{Z}_{i-1}, \bar{X}_{i-1}) = \sum_{i=2}^g H(\bar{Y}_{i-2}|\bar{Z}_{i-1}, \bar{X}_{i-2}). \quad (16)$$

Eq. (12) is by dislocation subtraction.  $\square$

The denominator in the definition  $\sum_{i=1}^g H(Y_i|\bar{X}_i, \bar{Y}_{i-1})$  represents the information quantity that is acquired in the search process. If function evaluations are independent, then  $H(Y_i|\bar{X}_i, \bar{Y}_{i-1}) = H(Y_i|X_i, \bar{X}_{i-1}, \bar{Y}_{i-1}) = H(Y_i|X_i)$ . While the numerator is more obscure. Actually it represents the quantity of the information about the objective function which is utilized by the algorithm (or in other words, the minimal information quantity that is needed to run the algorithm). Firstly,  $\sum_{i=1}^g H(Z_i|\bar{X}_{i-1}, \bar{Z}_{i-1}) = \sum_{i=1}^g H(Z_i|\bar{X}_{i-1}, \bar{Z}_{i-1}) - \sum_{i=1}^g H(Z_i|\bar{X}_{i-1}, \bar{Z}_{i-1}, \bar{Y}_{i-1})$  is similar to the concept of information gain in classification problems [10], which indicates the contribution of the information of  $Y$  to the algorithm. Secondly, the uncertainty of  $X$  and  $Z$  only lies in two aspects: the random sampling step and the lack of the information from  $Y$ . Thus  $H(X, Z) - \sum_{i=1}^g H(X_i|Z_i)$  can be regarded as the objective function's information that is utilized by the algorithm. And in fact, it is equal to the numerator. Thirdly, the numerator equals the denominator minus  $H(Y_g|\bar{X}_g, \bar{Y}_{g-1}) + H(\bar{Y}_{g-1}|Z, \bar{X}_{g-1})$  which can be seen as the wasted information of  $Y$ , because 1) the evaluation values in the last iteration  $Y_g$  cannot be utilized and 2) the information of previous evaluation values  $\bar{Y}_{g-1}$  is fully utilized only if  $H(\bar{Y}_{g-1}|Z, \bar{X}_{g-1}) = 0$ , i.e.,  $\bar{Y}_{g-1}$  can be reconstructed with  $Z$  given  $\bar{X}_{g-1}$ .

### 3 IURs of Heuristic Optimization Algorithms

In order to calculate the IURs of algorithms, we further assume  $f(x) \in \mathcal{Y}$  is identically and independently distributed (i.i.d). In most cases, it is unwise to calculate the IUR by definition. To calculate the denominator is quite straightforward under the above assumption, which equals the number of evaluations times  $H(f(x))$ . For example, if there are 100 cities in a travelling salesman problem [11] and  $f(x)$  obey uniform distribution, then  $|\mathcal{Y}| = 100!$ ,  $H(f(x)) = \log 100!$ . While on the other hand, to directly calculate the numerator is difficult and unnecessary. In each iteration, the output  $\mathcal{A}(D_{i-1})$  is a certain distribution, which is usually determined by some parameters in the algorithm. In fact we can certainly find (or construct) the set of intermediate parameters  $M_i$  such

that 1) there is a bijection from  $M_i$  to  $Z_i$  given  $\bar{X}_{i-1}$  and 2)  $M_i$  is determined only by  $\bar{Y}_{i-1}$  (otherwise  $H(Z_i|\bar{X}_{i-1}, \bar{Y}_{i-1}) > 0$ ), then

$$\sum_{i=2}^g H(Z_i|\bar{X}_{i-1}, \bar{Z}_{i-1}) = \sum_{i=2}^g H(M_i|\bar{M}_{i-1}) = H(M). \quad (17)$$

We only have to know the information quantity that is required to determine these intermediate parameters.

In the following, we investigate the IURs of several heuristics to show the procedure of calculating the IUR. Although these algorithms are designed for continuous (domain) optimization, the IURs of any kind of (discrete, combinatorial, dynamic, multi-objective) optimization algorithms can be calculated in the same way as long as there are a domain and a codomain. Without loss of generality, the following algorithms are all minimization algorithms, that is, they all intend to find the solution with the minimal evaluation value in the search space.

### 3.1 Random Search Algorithms

#### 3.1.1 Monte Carlo

The Monte Carlo (MC) method is often considered as a baseline for optimization algorithms. It is not a heuristic algorithm and usually fails to find acceptable solutions. If the maximal evaluation number is  $m$ , MC just uniformly randomly sample  $m$  solutions from  $\mathcal{X}$ .

MC does not utilize any information about the objective function because  $Z$  is fixed.

**Proposition 1.**

$$IUR_{MC} = 0. \quad (18)$$

#### 3.1.2 Luus-Jaakola

Luus-Jaakola (LJ) [12] is a heuristic algorithm based on MC. In each iteration, the algorithm generates a new individual  $y$  with the uniform distribution within a hypercube whose center is the position of the current individual  $x$ . If  $f(y) < f(x)$ ,  $x$  is replaced by  $y$ ; otherwise, the radius of the hypercube is multiplied by a parameter  $\gamma < 1$ .

The output of LJ in each iteration is the uniform distribution within the hypercube, which is determined by the position  $x$  and the radius. They are both controlled by the comparison result, i.e.,  $I(f(y) < f(x))$ .  $f(y)$  is i.i.d, but  $f(x)$  is the best in the history. Thus,  $H(M_i|\bar{M}_{i-1}) = H(I(f(y) < f(x))|\bar{M}_{i-1}) = \pi(i-1)$ .

**Proposition 2.**

$$IUR_{LJ}(g) = \frac{\sum_{i=1}^{g-1} \pi(i)}{gH(f(x))}. \quad (19)$$

## 3.2 Evolution Strategies

### 3.2.1 $(\mu, \lambda)$ -Evolution Strategy

$(\mu, \lambda)$ -evolution strategy (ES) [13] is an important heuristic algorithm in the family of evolution strategies. In each generation,  $\lambda$  new offspring are generated from  $\mu$  parents by crossover and mutation with normal distribution, and then the parents of a new generation are selected from these  $\lambda$  offspring. As a self-adaptive algorithm, the step size of the mutation is itself mutated along with the position of an individual.

The distribution for generating new offspring is determined by the  $\mu$  parents, namely the indexes of the best  $\mu$  of the  $\lambda$  individuals. Each set of  $\mu$  candidates has the same probability to be the best.  $H(M_i|\overline{M}_{i-1}) = H(M_i) = \log \binom{\lambda}{\mu}$ , where  $\binom{\lambda}{\mu} = \frac{\lambda!}{\mu!(\lambda-\mu)!}$ .

**Proposition 3.**

$$IUR_{(\mu, \lambda)\text{-ES}}(g) = \frac{(g-1) \log \binom{\lambda}{\mu}}{g\lambda H(f(x))}. \quad (20)$$

### 3.2.2 Covariance Matrix Adaptation Evolution Strategy

In order to more adaptively control the mutation parameters in  $(\mu, \lambda)$ -ES, a covariance matrix adaptation evolution strategy (CMA-ES) was proposed [14]. CMA-ES is a very complicated estimation of distribution algorithm [15], which adopts several different mechanisms to adapt the mean, the covariance matrix and the step size of the mutation operation. It is very efficient on benchmark functions especially when restart mechanisms are adopted. CMA-ES cannot be introduced here in detail. We refer interested readers to an elementary tutorial: [14].

Given  $\overline{X}_{i-1}$ , the mean, the covariance matrix and the step size of the distribution is determined by the indexes and the rankings of the best  $\mu$  individuals in each iteration in history.  $H(M_i|\overline{M}_{i-1}) = \log \frac{\lambda!}{(\lambda-\mu)!}$ .

**Proposition 4.**

$$IUR_{CMA\text{-ES}}(g) = \frac{(g-1) \log \frac{\lambda!}{(\lambda-\mu)!}}{g\lambda H(f(x))}. \quad (21)$$

Compared with  $(\mu, \lambda)$ -ES, it is obvious that  $IUR_{CMA\text{-ES}} \geq IUR_{(\mu, \lambda)\text{-ES}}$ , because not only the indexes of the  $\mu$  best individuals, but also their rankings are used in CMA-ES (to calculate their weights, for example). By utilizing the information of the solutions more thoroughly, CMA-ES is able to obtain more accurate knowledge of the objective function and search more efficiently.

The IURs of Particle Swarm algorithms [16, 17] and Differential Evolution algorithms [2, 5] are also investigated, shown in the appendix. If readers are interested in the IURs of other algorithms, we encourage you to conduct an investigation on your own which can be usually done with limited effort.



## 4 IUR versus Performance

The IUR is an intrinsic property of a heuristic algorithm, but the performance is not. Besides the algorithm itself, the performance of a heuristic also depends on the termination criterion, the way to measure the performance, and most importantly the distribution of the objective functions. A well-known fact about performance is that no algorithm outperforms another when there is no prior distribution [18], which is quite counter-experience. The objective functions in the real world usually subject to a certain underlying distribution. Although it is usually very difficult to precisely describe this distribution, we know that it has a much smaller information entropy than the uniform distribution and hence there is a free lunch [19, 20, 21]. In this case, the objective function (and resultantly its optimal point) can be identified with limited information quantity (the entropy of the distribution).

Reconsider the setting of the no free lunch (NFL) theorem from the perspective of information utilization. Assume  $|\mathcal{X}| = m$  and  $|\mathcal{Y}| = n$ . Under the setting of NFL (no prior distribution), the total uncertainty of the objective function is  $\log n^m = m \log n$ . In each evaluation, the information acquired is  $\log n$ . Therefore, no algorithm is able to certainly find the optimal point of the objective function within less than  $m$  times of evaluation even if all acquired information is thoroughly utilized. In this case, enumeration is the best algorithm [22]. On the contrary, if we already know the objective function is a sphere function, which is determined only by its center, then the required information quantity is  $\log m$ , and the least required number of evaluations is (more than)  $\log m / \log n = \log_n m$ . Suppose the dimensionality of the search space is  $d$ , then  $n$  is  $O(m^{\frac{1}{d}})$ ,  $\log_n m$  is  $O(d)$ , which is usually acceptable. If information is fully utilized ( $IUR \approx 1$ ), the exact number is  $d + 1$  [20]. While for algorithms with smaller IURs, more evaluations are needed. For example, if the IUR of another algorithm is half of the best algorithm (with half of the acquired information wasted), then at least about  $2d + 2$  evaluations are needed.

How much information is utilized by the algorithm per each evaluation determines the lower bound of the required evaluation number to locate the optimal point. **In this sense, IUR determines the upper bound of an algorithm's performance.** That is, algorithms with larger IURs have greater potential. However, the actual performance also depends on the manner of information utilization and how it accords with the underlying distribution of the objective function. For instance, one can easily design an algorithm with the same IUR as CMA-ES but does not work.

In the following, we will give empirical evidences on the correlation between IUR and performance. The preconditions of the experiments include: 1) the algorithms we investigate here are reasonably designed to optimize the objective functions from the underlying distribution; 2) the benchmark suite is large and comprehensive enough to represent the underlying distribution. The following conclusions may not hold for algorithms that are not reasonably designed or for a narrow or special range of objective functions. In other words, if the manner of information utilization does not accord with the underlying distribution of

objective functions, utilizing more information is not necessarily advantageous.

The theoretical correctness of IUR does not rely on these experimental results, but these examples may help readers understand how and to what extent IUR influences performance.

## 4.1 Different Parameters of the Same Algorithm

Sometimes for a certain optimization algorithm the IUR is influenced by only a few parameters. For these algorithms, we may adapt these parameters to show the correlation between the tendency of IUR and the tendency of performance.

### 4.1.1 $(\mu, \lambda)$ -ES

Intuitively, using  $\mu = \lambda$  is not a sensible option for  $(\mu, \lambda)$ -ES (commonly used  $\mu/\lambda$  values are in the range from  $1/7$  to  $1/2$  [23]) because it makes the selection operation invalid. Now we have a clearer explanation: the IUR of  $(\mu, \lambda)$ -ES is zero if  $\mu = \lambda$  (see Eq. (20)), i.e.,  $(\mu, \lambda)$ -ES does not use any heuristic information if  $\mu = \lambda$ .

Using a  $\mu$  around  $\frac{1}{2}\lambda$  may be a good choice for  $(\mu, \lambda)$ -ES because it leads to a large IUR. When  $\mu = \frac{1}{2}\lambda$ , the information used by  $(\mu, \lambda)$ -ES is the most. From the perspective of exploration and exploitation, we may come to a similar conclusion. If  $\mu$  is too small (elitism), the information of the population is only used to select the best few solutions, and resultantly the diversity of the population may suffer quickly. If  $\mu$  is too large (populism), the information of the population is only used to eliminate the worst few solutions, and resultantly the convergence speed may be too slow.

Different values of  $\mu/\lambda$  are evaluated on the CEC 2013 benchmark suite containing 28 different test functions (see Table 1) which are considered as black-box problems [24]. The meta parameter is set to  $\Delta\sigma = 0.5$ . The algorithm using each set of parameters is run 20 times independently for each function. The dimensionality is  $d = 5$ , and the maximal number of function evaluations is  $10000d$  for each run. For each fixed  $\lambda$ , the mean errors of 20 independent runs of each  $\mu/\lambda$  are ranked. The rankings are averaged over 28 functions, shown in Fig. 2.  $-\log\left(\frac{\lambda}{\mu}\right)/\lambda$  curves are translated along the vertical axis, also shown in Fig. 2.

According to the experimental results,  $\mu/\lambda$  around 0.5 is a good choice, which accord with our expectation. Moreover, the tendency of the performance (average ranking) curve is generally identical to that of the IUR ( $IUR_{(\mu,\lambda)-ES} \propto \log\left(\frac{\lambda}{\mu}\right)/\lambda$  when  $g$  is large). The experimental results indicate a positive correlation between the performance and the IUR: the parameter value with larger IUR is prone to perform better. Thus, the IUR can be used to guide the choice of parameters. After all, tuning the parameters by experiments is much more expensive than calculating the IURs.

Table 1: Test functions of CEC 2013 single objective optimization benchmark suite [24]

	No.	Name
Unimodal Functions	1	Sphere Function
	2	Rotated High Conditioned Elliptic Function
	3	Rotated Bent Cigar Function
	4	Rotated Discus Function
	5	Different Powers Function
Basic Multimodal Functions	6	Rotated Rosenbrock's Function
	7	Rotated Schaffers F7 Function
	8	Rotated Ackley's Function
	9	Rotated Weierstrass Function
	10	Rotated Griewank's Function
	11	Rastrigin's Function
	12	Rotated Rastrigin's Function
	13	Non-Continuous Rotated Rastrigin's Function
	14	Schwefel's Function
	15	Rotated Schwefel's Function
	16	Rotated Katsuura Function
	17	Lunacek Bi_Rastrigin Function
	18	Rotated Lunacek Bi_Rastrigin Function
	19	Expanded Griewank's plus Rosenbrock's Function
	20	Expanded Scaffer's F6 Function
Composition Functions	21	Composition Function 1 (Rotated)
	22	Composition Function 2 (Unrotated)
	23	Composition Function 3 (Rotated)
	24	Composition Function 4 (Rotated)
	25	Composition Function 5 (Rotated)
	26	Composition Function 6 (Rotated)
	27	Composition Function 7 (Rotated)
	28	Composition Function 8 (Rotated)

#### 4.1.2 CMA-ES

Different from  $(\mu, \lambda)$ -ES, CMA-ES adopts a rank-based weighted recombination instead of a selection operation, in which the rank information of the best  $\mu$  individuals is utilized.

On the one hand, the rank-based weighted recombination achieves the largest IUR when  $\mu = \lambda$  (see Eq. (21)). The larger  $\mu$  is, the more information is used (because the rank information of the rest  $\lambda - \mu$  individuals are wasted).

On the other hand, the rank-based weighted recombination also achieves the best performance when  $\mu = \lambda$  [25]. This is also an evidence on the correlation between IUR and performance. However, the optimal weighted recombination requires the use of negative weights, which is somehow not adopted in CMA-ES [14]. The manner of information utilization and other conditions (termination criterion, performance measure, etc.) should also be taken into consideration when parameters are chosen. Therefore using  $\mu = \lambda$  is probably not the best choice for CMA-ES even though it leads to a large IUR.

## 4.2 Algorithms in the Same Family

Usually different algorithms in the same family utilize information in similar manners, in which case we may compare their performances to show the correlation between IUR and performance. However, we need to be more cautious

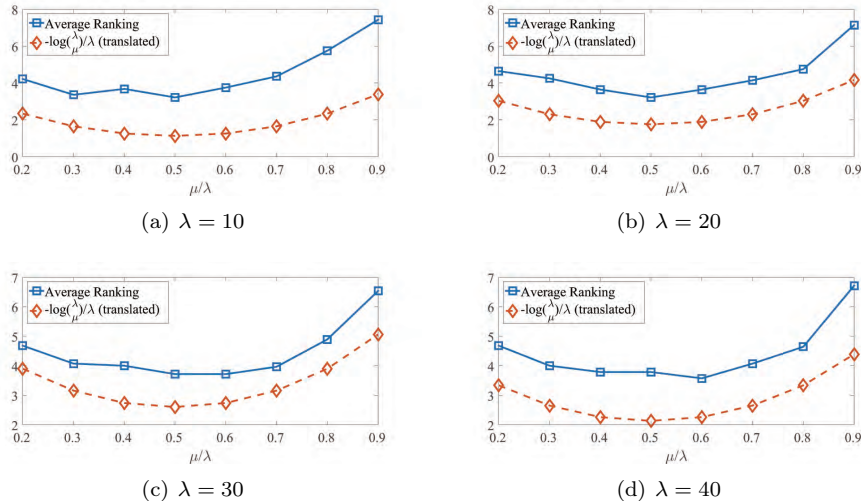


Figure 2: The average rankings and the  $-\log(\frac{\lambda}{\mu})/\lambda$  curves of each value of  $\lambda$ .

here because IUR is not the only factor as long as different algorithms in the same family do not utilize information in identical manners.

As shown in Section 3,  $IUR_{LJ} \geq IUR_{MC}$  and  $IUR_{CMA-ES} \geq IUR_{(\mu, \lambda)-ES}$ . LJ and CMA-ES are more finely designed compared with the previous algorithms since they are able to utilize more information of the objective function. Naturally we would expect that LJ outperforms MC and CMA-ES outperforms  $(\mu, \lambda)$ -ES.

The four algorithms are evaluated on the CEC 2013 benchmark suite. The parameter of LJ is set to  $\gamma = 0.99$ . The parameters of  $(\mu, \lambda)$ -ES are set to  $\lambda = 30, \mu = 15, \Delta\sigma = 0.5$ . The parameters of CMA-ES are set to suggested values [14] except that  $\sigma = 50$  because the radius of the search space is 100. The dimensionality is  $d = 5$ , and the maximal number of function evaluation is  $10000d$  for each run. Each algorithm is run 20 times independently for each function. Their mean errors are shown in Table 2. The best mean errors are highlighted. Their mean errors are ranked on each function, and the average rankings (AR.) over 28 functions are also shown in Table 2.

Pair-wise Wilcoxon rank sum tests are also conducted between MC and LJ and between  $(\mu, \lambda)$ -ES and CMA-ES. The  $p$  values are shown in the last two columns of Table 2. Significant results (with confidence level 95%) are underlined. The results of LJ are significantly better than MC on 14 functions, and significantly worse on only 10 functions. While the results of CMA-ES are significantly better than  $(\mu, \lambda)$ -ES on 14 functions, and significantly worse on only 3 functions. Generally speaking, the performance of LJ is better than MC and the performance of CMA-ES is better than  $(\mu, \lambda)$ -ES. These experimental results imply that the extent of information utilization may be an important

Table 2: Mean errors and average rankings of the four algorithms and  $p$  values

F.	MC	LJ	$(\mu, \lambda)$ -ES	CMA-ES	MC vs. LJ	$(\mu, \lambda)$ -ES vs. CMA-ES
1	2.18E+02	<b>0.00E+00</b>	5.91E-12	<b>0.00E+00</b>	<u>8.01E-09</u>	<u>4.01E-02</u>
2	4.25E+05	<b>0.00E+00</b>	3.49E+05	<b>0.00E+00</b>	<u>8.01E-09</u>	<u>8.01E-09</u>
3	8.02E+07	<b>0.00E+00</b>	2.18E+07	<b>0.00E+00</b>	<u>1.13E-08</u>	<u>1.13E-08</u>
4	4.23E+03	<b>0.00E+00</b>	2.20E+04	<b>0.00E+00</b>	<u>1.13E-08</u>	<u>1.13E-08</u>
5	8.00E+01	6.79E+01	1.95E-05	<b>0.00E+00</b>	<u>1.33E-02</u>	<u>1.90E-04</u>
6	9.94E+00	2.51E+01	<b>2.46E+00</b>	7.86E-01	<u>4.17E-05</u>	<u>8.15E-06</u>
7	2.02E+01	7.10E+01	1.66E+01	<b>5.66E+00</b>	1.99E-01	<u>2.56E-03</u>
8	<b>1.83E+01</b>	2.01E+01	2.03E+01	2.10E+01	<u>3.42E-07</u>	<u>1.61E-04</u>
9	2.53E+00	1.67E+00	2.37E+00	<b>1.08E+00</b>	<u>1.48E-03</u>	<u>1.63E-03</u>
10	2.30E+01	1.78E+00	1.30E+01	<b>4.16E-02</b>	<u>6.80E-08</u>	<u>1.23E-07</u>
11	2.22E+01	1.40E+01	6.67E+00	<b>6.57E+00</b>	<u>3.04E-04</u>	8.17E-01
12	2.10E+01	1.33E+01	1.20E+01	<b>7.36E+00</b>	<u>1.12E-03</u>	<u>2.04E-02</u>
13	2.21E+01	1.90E+01	1.87E+01	<b>1.28E+01</b>	1.20E-01	5.98E-01
14	3.78E+02	7.53E+02	<b>1.35E+02</b>	4.61E+02	<u>1.10E-05</u>	<u>7.41E-05</u>
15	<b>3.84E+02</b>	6.85E+02	5.27E+02	4.52E+02	<u>3.99E-06</u>	1.81E-01
16	7.43E-01	<b>5.34E-01</b>	8.27E-01	1.49E+00	<u>1.93E-02</u>	1.11E-01
17	3.25E+01	2.23E+01	<b>9.87E+00</b>	1.07E+01	<u>1.78E-03</u>	3.65E-01
18	3.43E+01	1.82E+01	<b>1.01E+01</b>	1.01E+01	<u>2.60E-05</u>	9.89E-01
19	4.08E+00	7.21E-01	5.45E-01	<b>4.82E-01</b>	<u>9.17E-08</u>	9.46E-01
20	<b>1.23E+00</b>	1.85E+00	2.50E+00	1.92E+00	<u>1.10E-05</u>	<u>6.97E-06</u>
21	3.23E+02	3.05E+02	<b>2.55E+02</b>	2.80E+02	<u>1.94E-02</u>	9.89E-01
22	5.91E+02	7.91E+02	<b>4.01E+02</b>	7.20E+02	<u>2.56E-03</u>	<u>5.63E-04</u>
23	<b>6.04E+02</b>	8.33E+02	7.01E+02	6.08E+02	<u>8.29E-05</u>	3.37E-01
24	<b>1.26E+02</b>	2.04E+02	1.99E+02	1.76E+02	<u>6.80E-08</u>	<u>4.60E-04</u>
25	<b>1.27E+02</b>	1.96E+02	1.98E+02	1.81E+02	<u>1.60E-05</u>	<u>7.71E-03</u>
26	<b>1.01E+02</b>	2.38E+02	1.67E+02	1.98E+02	<u>1.43E-07</u>	7.76E-01
27	3.57E+02	3.52E+02	3.65E+02	<b>3.27E+02</b>	4.25E-01	<u>2.47E-04</u>
28	3.05E+02	<b>3.00E+02</b>	3.25E+02	3.15E+02	8.59E-01	<u>2.03E-01</u>
AR.	2.82	2.68	2.43	<b>1.82</b>	14 : 10	14 : 3

factor in the performance.

The algorithms in the same family utilize information in similar but different manners. In this case, the influence of IUR on the performance is crucial, but sometimes not deterministic. For example, in CMA-ES, there are several different mechanisms proposed to improve the performance. The improvement in IUR does not reflect all of them. The improvement related to IUR is the rank-based weighted recombination. It has significant impact on performance [25, 26]. While other mechanisms such as adapting the covariance matrix and the step size are not related to IUR but also very important. These mechanisms are introduced as different information utilization manners, which help the algorithm to better fit the underlying distribution of objective functions. Similar comparisons can be made between PSO and SPSO and between DE and JADE (see appendix).

However, after all, an algorithm cannot perform very well if little information is used. Hence, just like LJ, CMA-ES, SPSO and JADE, the tendency of elevating the IUR is quite clear in various families of heuristics. Many mechanisms have been proposed to better preserve historical information for further utilization [27, 28, 5]. Many general methods (adaptive parameter control [29], estimation of distribution [15], fitness approximation [30], Bayesian approaches [31], Gaussian process models [32], hyper-heuristic [33]) have been proposed to elevate the IURs of heuristics. Not to mention these numerous specified mech-

anisms. In summary, the IUR provides an important and sensible perspective on the developments in this field.

### 4.3 Algorithms in Different Families

The correlation between the IUR and the performance of the algorithms in different families (such as LJ and  $(\mu, \lambda)$ -ES) can be vaguer because the manners of information utilization are different, though the above experimental results accord with our expectation ( $(\mu, \lambda)$ -ES performs better than LJ and  $IUR_{LJ} \leq IUR_{(\mu, \lambda)\text{-ES}}$  unless  $\mu = \lambda$ ). If algorithms utilize information in extremely different manners, the IUR may not be the deterministic factor. There are infinite manners to utilize information. It is difficult to judge which manner is better. Whether a manner is good or not depends on how it fits the underlying distribution of the objective functions, which is difficult to describe. A well designed algorithm with low IUR may outperform a poorly designed algorithm with high IUR because it utilizes information more efficiently and fits the underlying distribution better. Nonetheless, certainly the extent of information utilization is still of importance in this case because 1) the algorithms with larger IURs have greater potential 2) the IUR of the "best" algorithm (if any) must be very close to one and 3) an algorithm that uses little information cannot be a good algorithm.

The exact correlation between the IUR and the performance requires much more theoretical works on investigating the manners of information utilization and how they fit the underlying distributions, which are very difficult but not impossible.

## 5 Upper Bound for Comparison-based Algorithms

Above examples have covered several approaches of information utilization in heuristic optimization algorithms. But the IURs of these algorithms are all not high because they are comparison-based algorithms, in which only the rank information is utilized.

**Theorem 2** (Upper bound for comparison-based algorithms). *If the maximal number of evaluations is  $m$ ,  $y = f(x)$  are i.i.d, and algorithm  $\mathcal{A}$  is a comparison-based optimization algorithm,*

$$IUR_{\mathcal{A}} \leq \frac{\log m}{H(f(x))}. \quad (22)$$

*Proof.* Suppose in a certain run, the actual evaluation number is  $m' \leq m$ . In this case,  $M$  is drawn from a set with cardinal number at most  $m'!$  (with  $m'$  individuals all sorted), then the maximal information quantity is  $H(M) \leq \log m'!$  for a comparison-based algorithm. Thus  $IUR_{\mathcal{A}} \leq \frac{\log m'!}{m'H(f(x))}$ . Note that the right hand side is a monotonically increasing function of  $m'$ , and  $\frac{\log m'!}{m'H(f(x))} \leq \frac{\log m}{mH(f(x))} \leq \frac{\log m}{H(f(x))}$ .  $\square$

Suppose  $|\mathcal{Y}| = n$  and  $f(x)$  obey uniform distribution, then  $\frac{\log m}{H(f(x))} = \log_n m$ . Typically  $m \ll n$ , thus this upper bound is quite low. Most iterative algorithms do not allow the information in past iterations (because it requires a lot of memory space to do so), in which case the upper bound becomes  $\frac{\log \lambda}{H(f(x))}$  where  $\lambda$  is the evaluation number in each generation. The IUR of CMA-ES is able to approach this bound when  $\mu = \lambda$ . That is, CMA-ES has almost the largest IUR in comparison-based algorithms without historical information.

There exist algorithms which use exact evaluation values in the searching process, such as genetic algorithm [34], ant colony optimization [35], estimation of distribution algorithms [15], invasive weed optimization [36], artificial bee colony [37], fireworks algorithm [38], etc. They can achieve higher IURs, even close to 1, because the cardinal number of the set from which  $M$  is drawn can be up to  $n^m$ . These algorithms have greater potential than comparison-based algorithms and can outperform them if well designed.

## 6 Conclusion

It is natural and often effective to utilize more heuristic information in optimization algorithms, which has been widely realized. However, there was no metric to reflect the extent of information utilization. In this paper, a metric called the information utilization ratio (IUR) is defined as the ratio of the utilized information quantity over the acquired information quantity. IUR can be an index to reflect how finely and advanced an algorithm is designed. IUR proves to be well defined. Several examples are given to demonstrate the procedure of calculating IURs. Generally speaking, the IUR determines the upper bound of the performance of an optimization algorithm. To further indicate the importance of this metric, several experiments are conducted to show the correlation between the IUR and the performance. The experimental results imply that 1) for a certain algorithm, the parameter value with larger IUR has advantage; 2) for algorithms in the same family, the one with larger IUR is prone to be more efficient; 3) for algorithms in different families, the IUR is also an important factor. We also give the IUR's upper bound for comparison-based algorithms.

The IUR can be used to guide the choice of parameters, guide the design of new algorithms and guide the improvement of existing algorithms. For example, if you are inventing a new algorithm, or adapting an existing one, it is promising to include mechanisms that can enhance the information utilization in your algorithm. If you want to know which one among several algorithms is more likely efficient before you use them, it would be quite informative to compare their IURs to show which one is better designed and has greater potential.

Most works in the field of heuristic search or optimization focus on inventing new mechanisms or tricks, while few have considered the potential driver behind these works. We consider this work as a fundamental theory, which is surprisingly not easy. Hopefully the definition of IUR will lead to a more systematic manner of research about how mechanisms should be designed and how information should be utilized.

Extending this metric to other fields in artificial intelligence such as classification and time series prediction may be an interesting future work.



## .1 Particle Swarm Algorithms

### .1.1 Particle Swarm Optimization

Particle swarm optimization (PSO) [16] is one of the most famous swarm and heuristic algorithms which is quite simple but surprisingly efficient in numerical optimization. In PSO, a fixed number ( $s$ ) of particles moves in the search space to find the optimal solutions. The position of a particle is updated as follows. In generation  $g$ , for each particle  $i$  and each dimension  $j$ ,

$$v_{ij}(g+1) \leftarrow v_{ij}(g) + \phi_1 r_{1,ij}(pbest_{ij}(g) - x_{ij}(g)) + \phi_2 r_{2,ij}(gbest_j(g) - x_{ij}(g)), \quad (23)$$

$$x_{ij}(g+1) \leftarrow x_{ij}(g) + v_{ij}(g+1), \quad (24)$$

where  $\phi_1$  and  $\phi_2$  are constant coefficients,  $r_1$  and  $r_2$  are random numbers,  $pbest$  is the best position in history found by this particle and  $gbest$  is the best position found by the entire swarm.

The output distribution in each generation is determined by  $I(f(x_i(g)) < f(pbest_i(g-1)))$  and  $\arg \min_i f(pbest_i(g))$ . Although it is difficult to calculate  $H(M)$ , we have the lower and upper bounds:

$$s \sum_{i=1}^{g-1} \pi(i) \leq H(M) \leq \sum_{i=2}^g H(M_i) \leq (g-1) \log s + s \sum_{i=1}^{g-1} \pi(i). \quad (25)$$

**Proposition 5.**

$$\frac{s \sum_{i=1}^{g-1} \pi(i)}{sgH(f(x))} \leq IUR_{PSO}(g) \leq \frac{(g-1) \log s + s \sum_{i=1}^{g-1} \pi(i)}{sgH(f(x))}. \quad (26)$$

### .1.2 Standard Particle Swarm Optimization

After years of development, many improvements and variants are proposed for PSO. In order to construct a common ground for further researches, a standard particle swarm optimization (SPSO) was defined [17]. Compared with original PSO, there are two main modifications: the local ring topology and the constricted update rule. The constricted update rule uses a new coefficient derived from  $\phi_1$  and  $\phi_2$  to constrict the velocity to guarantee convergence. In the local ring topology, the  $gbest$  in the velocity update equation is replaced with a  $lbest$ , which is the best position among this individual and its two neighbourhoods on the ring.

For each group (consisting of three particles), information with quantity at most  $\log 3$  is needed to decide  $lbest$ .

**Proposition 6.**

$$\frac{s \sum_{i=1}^{g-1} \pi(i)}{sgH(f(x))} \leq IUR_{SPSO}(g) \leq \frac{s(g-1) \log 3 + s \sum_{i=1}^{g-1} \pi(i)}{sgH(f(x))}. \quad (27)$$

Usually  $IUR_{PSO} \leq IUR_{SPSO}$  though their exact values are difficult to derive. It turns out that the information utilization ratio of the local model is larger than the global model because in local topology the particles interact with each other more frequently.

According to experimental results, SPSO significantly outperform PSO on a large range of test functions. [17]

## .2 Differential Evolution Algorithms

### .2.1 Differential Evolution

Differential evolution (DE) [2] is a powerful heuristic algorithm for numerical optimization. The number of individuals in DE is also fixed. The mutation is conducted as below (take DE/rand/1 as an example). For each  $x$  in the population, generate

$$z = x_{r1} + F(x_{r2} - x_{r3}), \quad (28)$$

where  $r1, r2$  and  $r3$  are random indexes and  $F$  is a constant coefficient. Then a crossover is conducted between  $z$  and  $x$  to generate a new candidate  $y$ , where there is a parameter  $CR$  to control the probability that a dimension of  $y$  is identical to that of  $z$ . If  $f(y) < f(x)$ ,  $x$  is replaced with  $y$ , otherwise,  $x$  is kept.

In DE, the distribution of generating new offspring is determined by  $I(f(y) < f(x))$  of each individual. So the IUR of DE is equal to that of LJ with the same  $g$ . However, they would be different with the same number of evaluation times.

**Proposition 7.**

$$IUR_{DE}(g) = \frac{s \sum_{i=1}^{g-1} \pi(i)}{sgH(f(x))}. \quad (29)$$

IURs of some other DE variants are given in Table 3.

Table 3: IURs of other DE variants

	IUR
DE/best/1	$= IUR_{PSO}$
DE/current-to-best/1	$= IUR_{PSO}$
DE/rand/2	$= IUR_{DE}$
DE/best/2	$= IUR_{PSO}$

### .2.2 JADE

JADE [5] is an important development of DE. There are three main adaptations proposed in JADE:

1. A DE/current-to- $p$ best/1 mutation strategy. In JADE,

$$z_i = x_i + F_i(x_{best}^p - x_i) + F_i(x_{r1} - x_{r2}), \quad (30)$$

where  $x_{best}^p$  is a randomly chosen individual from the 100 $p$ % best individuals.

2. An optional external archive.
3. Adaptive mutation parameters.

External archive is a useful tool to improve information utilization. However, in JADE these individuals are just randomly chosen and randomly removed from the archive, where no information of the objective function is used. Compared to DE, JADE elevates IUR after all because the indexes of the best 100p% individuals are used. Note that the output distribution is determined only when all indexes of the best 100p% individuals are given.

**Proposition 8.**

$$\frac{s \sum_{i=1}^{g-1} \pi(i)}{sgH(f(x))} \leq IUR_{JADE}(g) \leq \frac{(g-1) \log \binom{s}{ps} + s \sum_{i=1}^{g-1} \pi(i)}{sgH(f(x))}. \quad (31)$$

According to experimental results, JADE significantly outperform DE on a large range of test functions. [5]

## References

- [1] Mandavilli Srinivas and Lalit M Patnaik. Genetic algorithms: A survey. *Computer*, 27(6):17–26, 1994.
- [2] Rainer Storn and Kenneth Price. Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces. *Journal of global optimization*, 11(4):341–359, 1997.
- [3] Marco Dorigo, Mauro Birattari, and Thomas Stützle. Ant colony optimization. *Computational Intelligence Magazine, IEEE*, 1(4):28–39, 2006.
- [4] Georges R Harik, Fernando G Lobo, and David E Goldberg. The compact genetic algorithm. *Evolutionary Computation, IEEE Transactions on*, 3(4):287–297, 1999.
- [5] Jingqiao Zhang and Arthur C Sanderson. JADE: adaptive differential evolution with optional external archive. *Evolutionary Computation, IEEE Transactions on*, 13(5):945–958, 2009.
- [6] J. Li, S. Zheng, and Y. Tan. The effect of information utilization: Introducing a novel guiding spark in the fireworks algorithm. *IEEE Transactions on Evolutionary Computation*, PP(99):1–1, 2016.
- [7] Donald R. Jones. A taxonomy of global optimization methods based on response surfaces. *Journal of Global Optimization*, 21(4):345–383, 2001.
- [8] Jun He and Guangming Lin. Average convergence rate of evolutionary algorithms. *IEEE Transactions on Evolutionary Computation*, 20(2):1, 2015.

- [9] C.E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27(4):623–656, Oct 1948.
- [10] J. R. Quinlan. Induction of decision trees. *Machine Learning*, 1(1):81–106, 1986.
- [11] Eugene L Lawler. The traveling salesman problem: a guided tour of combinatorial optimization. *WILEY-INTERSCIENCE SERIES IN DISCRETE MATHEMATICS*, 1985.
- [12] Rein Luus and THI Jaakola. Optimization by direct search and systematic reduction of the size of search region. *AIChE Journal*, 19(4):760–766, 1973.
- [13] Thomas Bäck, Frank Hoffmeister, and Hans-Paul Schwefel. A survey of evolution strategies. In *Proceedings of the Fourth International Conference on Genetic Algorithms*, 1991.
- [14] Nikolaus Hansen. The CMA evolution strategy: A tutorial. *Vu le*, 29, 2005.
- [15] Pedro Larranaga. A review on estimation of distribution algorithms. In *Estimation of distribution algorithms*, pages 57–100. Springer, 2002.
- [16] Russell Eberhart and James Kennedy. A new optimizer using particle swarm theory. In *Micro Machine and Human Science, 1995. MHS'95., Proceedings of the Sixth International Symposium on*, pages 39–43. IEEE, 1995.
- [17] Daniel Bratton and James Kennedy. Defining a standard for particle swarm optimization. In *Swarm Intelligence Symposium, 2007. SIS 2007. IEEE*, pages 120–127. IEEE, 2007.
- [18] David H Wolpert and William G Macready. No free lunch theorems for optimization. *Evolutionary Computation, IEEE Transactions on*, 1(1):67–82, 1997.
- [19] Matthew J. Streeter. *Two Broad Classes of Functions for Which a No Free Lunch Result Does Not Hold*. Springer Berlin Heidelberg, 2003.
- [20] Anne Auger and Olivier Teytaud. Continuous lunches are free plus the design of optimal optimization algorithms. *Algorithmica*, 57(1):121–146, 2010.
- [21] Tom Everitt, Tor Lattimore, and Marcus Hutter. Free lunch for optimisation under the universal distribution. In *Evolutionary Computation (CEC), 2014 IEEE Congress on*, pages 167–174. IEEE, 2014.
- [22] T. M. English. Some information theoretic results on evolutionary optimization. In *Evolutionary Computation, Cec 99, Proceedings of the Congress on*, 1999.

- [23] H. Beyer. Evolution strategies. *Scholarpedia*, 2(8):1965, 2007. revision #130731.
- [24] JJ Liang, BY Qu, PN Suganthan, and Alfredo G Hernández-Díaz. Problem definitions and evaluation criteria for the CEC 2013 special session on real-parameter optimization, 2013.
- [25] Dirk V. Arnold. Optimal weighted recombination. In *Foundations of Genetic Algorithms, International Workshop, Foga 2005, Aizu-Wakamatsu City, Japan, January 5-9, 2005, Revised Selected Papers*, pages 215–237, 2005.
- [26] Nikolaus Hansen and Stefan Kern. *Evaluating the CMA Evolution Strategy on Multimodal Test Functions*. Springer Berlin Heidelberg, 2004.
- [27] Kalyanmoy Deb, Samir Agrawal, Amrit Pratap, and T. Meyarivan. *A Fast Elitist Non-dominated Sorting Genetic Algorithm for Multi-objective Optimization: NSGA-II*. Springer Berlin Heidelberg, 2000.
- [28] Margarita Reyes-Sierra and Carlos A. Coello Coello. Multi-objective particle swarm optimizers: A survey of the state-of-the-art. *International Journal of Computational Intelligence Research*, 2(3):287–308, 2006.
- [29] Agoston Endre Eiben, Robert Hinterding, and Zbigniew Michalewicz. Parameter control in evolutionary algorithms. *Evolutionary Computation, IEEE Transactions on*, 3(2):124–141, 1999.
- [30] Yaochu Jin. A comprehensive survey of fitness approximation in evolutionary computation. *Soft computing*, 9(1):3–12, 2005.
- [31] Martin Pelikan. Bayesian optimization algorithm. In *Hierarchical Bayesian optimization algorithm*, pages 31–48. Springer, 2005.
- [32] Dirk Büche, Nicol N Schraudolph, and Petros Koumoutsakos. Accelerating evolutionary algorithms with gaussian process fitness function models. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, 35(2):183–194, 2005.
- [33] Edmund K Burke, Matthew Hyde, Graham Kendall, Gabriela Ochoa, Ender Özcan, and John R Woodward. A classification of hyper-heuristic approaches. In *Handbook of metaheuristics*, pages 449–468. Springer, 2010.
- [34] John H Holland. *Adaptation in natural and artificial systems: An introductory analysis with applications to biology, control, and artificial intelligence*. 1975.
- [35] Marco Dorigo, Vittorio Maniezzo, and Alberto Coloni. Ant system: optimization by a colony of cooperating agents. *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, 26(1):29–41, 1996.

- [36] Ali Reza Mehrabian and Caro Lucas. A novel numerical optimization algorithm inspired from weed colonization. *Ecological informatics*, 1(4):355–366, 2006.
- [37] Dervis Karaboga and Bahriye Basturk. A powerful and efficient algorithm for numerical function optimization: artificial bee colony (abc) algorithm. *Journal of global optimization*, 39(3):459–471, 2007.
- [38] Ying Tan and Yuanchun Zhu. Fireworks algorithm for optimization. In *Advances in Swarm Intelligence*, pages 355–364. Springer, 2010.