
Optimization of Location Allocation of Web Service using non-dominated sorting algorithm(NSGA-II)

Boxiong Tan

tanboxi@ecs.vuw.ac.nz

Department of Science and Engineering, Victoria of Wellington, Wellington, Zip, New Zealand

Hui Ma

hui.ma@ecs.vuw.ac.nz

Department of Science and Engineering, Victoria of Wellington, Wellington, Zip, New Zealand

Abstract

The abstract goes here. It should be about 200 words and give the reader a summary of the main contributions of the paper. Remember that readers may decide to read or not to read your paper based on what is in the abstract. The abstract never contains references.

Keywords

multi-objective, NSGA-II, genetic algorithms, evolutionary programming

1 Introduction

Web Services are considered as self-contained, self-describing, modular applications that can be published, located, and invoked across the Web Ran (2003). In recent years, web services technology is becoming increasingly popular because the convenience, low cost and capacity to be composed into high-level business processes Aboolian et al. (2009). Because of modularization and open interfaces, Web services facilitate the development of highly customizable and adaptable applications to meet business demand. Furthermore, Web services offer a convenient registration, search and discovery system. e.g. Universal Description, Discovery and Integration (UDDI).

With the ever increasing number of functional similar web services being available on the Internet, the web service providers (WSPs) are trying to improve the quality of service (QoS) to become competitive in the market. QoS also known as non-functional requirements to web services, is the degree to which a service meets specified requirements or user needs Zhou and Niemela (2006), such as response time, security and availability. Among numerous QoS measurements, service response time is a critical factor for many real-time services, e.g. traffic service or finance service. Service response time has two components: transmission time (variable with message size) and network latency Johansson (2000). Study Jamin et al. (2001) has shown that network latency is a significant component of service response delay. Ignoring network latency will underestimate response time by more than 80 percent. Since network latency is related to network topology as well as physical distance Huffaker et al. (2002). The network latency could also vary with the network topology changes. The only way to reduce the network latency is move the service to a location where has smaller network

latency to the user center. Hence, the WSPs need to consider which physical location to deploy their services so that it could minimize the cost as well as ensure the QoS.

The Web service location-allocation problem is essentially a multi-objective optimization problem Caramia (2008). Because of the confliction between service quality and deployment cost. Ideally, WSP could deploy their services to each user center in order to provide the best quality. That is, the more services deployed, the better the quality and the higher cost. This problem is considered as an NP-hard due to the fact that the combinatorial explosion of the search space Vanrompay et al. (2008).

Very few researches Aboolian et al. (2009); Sun and Koehler (2006) study this problem. Both studies try to solve the problem by integer linear programming techniques. In particular, Sun and Koehler (2006) has solved this problem by employ a modified Alternate Location-Allocation (ALA) algorithm. However, essentially ALA build upon integer programming and use greedy algorithm to optimize the result. As a result, the solution is easily stuck at local optima.

Evolutionary algorithms (EAs) have been used in solving multi objective optimization problems in recent years. EAs are ideal for solving multi objective optimization problems Desai et al. (2012), since EA works with a population of solutions, a simple EA can be extended to maintain a diverse set of solutions. With an emphasis for moving toward the true Pareto-optimal region, an EA can be used to find multiple Pareto-optimal solutions in one single simulation run Kanagarajan et al. (2008).

Hai Huang et al. (2014) proposed an enhanced genetic algorithm-based approach which make use of the integer scalarization technique to solve the multi-objective problem. The genetic algorithm (GA) is an EA that uses genetic operators to obtain optimal solutions without any assumptions about the search space. This algorithm solve the scalability problem in the dataset, however the integer scalarization technique Caramia (2008) has some disadvantages:

1. The decision maker needs to choose an appropriate weights for the objectives to retrieve a satisfactorily solution.
2. The algorithm does not produce an uniform spread of points on the Pareto curve. That is, all points are grouped in certain parts of the Pareto front.
3. Non-convex parts of the Pareto set cannot be reached by minimizing convex combinations of the object functions.

Evolutionary multi objective optimization (EMO) methodologies on the other hand, successfully avoid the above mentioned problems and demonstrated their usefulness in find a well-distributed set of near Pareto optimal solutions Aboolian et al. (2009). Non-dominated sorting GA (NSGA-II) Deb et al. (2002), Strength Pareto Evolutionary Algorithm 2 (SPEA-2) Deb et al. (2005) have become standard approaches. Some schemes based on particle swarm optimization approaches Elhossini et al. (2010); Huang et al. (2006) are also important. Among numerous EA approaches, NSGA-II is one of the most widely used methods for generating the Pareto frontier. NSGA-II implements elitism and uses a phenotype crowd comparison operator that keeps diversity without specifying any additional parameters Deb et al. (2006). In our approach, we apply a modified version of NSGA-II since the web service location-allocation is a discrete problem.

The main purpose of this research is to provide a framework that guides a WSP to locate its servers. We consider the problem faced by a WSP who has existing facilities but wishes to use the collected data to re-allocate their services in order to maximum

their profit. The WSP must decide on facility locations from a finite set of possible locations. In order to make the decision, the WSP must first analyze the data which were collected from current services. The collected data should includes the records of invocation from each unique IP address. Therefore, based on these data, the WSP could summarize several customer demand concentrated at n discrete nodes Aboolian et al. (2009), namely user centers. We assume the WSP has already done this step and list of user centers and candidate service deployment locations are given. In addition to decide which location to re-allocate the services, a dataset which contains the network latency between demand user center and candidate location are critical. The WSP could collect the data or use existed dataset Zheng et al. (2010); Zhang et al. (2011). Then, the service provider could use the algorithm which proposed by this paper, to select an optimal plan based on their funds. The algorithm will produce a near optimal solution which indicate the services deployment locations with a minimum cost and best service quality. The main objectives are:

- To model the web service location-allocation problem so that it can be tackled with NSGA-II
- To develop a modified NSGA-II approach for the web service location-allocation problem
- To evaluate our approach by comparing it to a GA approach which use integer scalarization technique.

In Section 2 we provide a formulation for our model. Section 3 develops a modified NSGA-II. We use these to study a number of hypotheses on the placement of WSPs. These solutions are compared to solutions from ALA. The experiments and their result are discussed in Section 4. We end with a discussion of our results and possible future research directions in Section 5.

2 Problem Description

2.1 Model formulation

The problem is to determine which facility locations that could maximus WSPs profit as well as ensure low network latency. Let $S = \{1, 2, \dots, s\}$ be the set of services. We assume that the demand for service is concentrated at i demand nodes $I = \{1, 2, \dots, i\}$. Let $J = \{1, 2, \dots, j\}$ be the set of j candidate facility locations. To model the service location-allocation problem we use four matrices: service network latency matrix L , service location matrix A , service invocation frequency matrix F and cost matrix C .

The server network latency matrix $L = [l_{ij}]$, is used to record network latency from user centers to candidate locations, where l_{ij} is a real number denotes the network latency from user center i to candidate location j . These data could be retrieved from implementing a network latency experiment or using existed datasets Zheng et al. (2010); Zhang et al. (2011).

$$L = \begin{matrix} & \begin{matrix} j_1 & j_2 & j_3 \end{matrix} \\ \begin{matrix} i_1 \\ i_2 \\ i_3 \end{matrix} & \begin{bmatrix} 0 & 5.776 & 6.984 \\ 5.776 & 0 & 2.035 \\ 0.984 & 1.135 & 2.3 \end{bmatrix} \end{matrix}$$

The service location matrix $A = [y_{sj}]$ represents the actual service location-allocation, where y_{sj} is a binary value (i.e., 1 or 0) shows whether a service s is deployed in

candidate location j or not. We use the service location matrix A as the representation of chromosome that evolve itself during the progress.

$$A = \begin{matrix} & j_1 & j_2 & j_3 \\ \begin{matrix} s_1 \\ s_2 \\ s_3 \end{matrix} & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \end{matrix}$$

The service invocation frequency matrix $F = [f_{is}]$, is used to record services invocation frequency from user centers, which f_{is} is an integer that indicate the number of invocation in a period of time from user center i to service s . e.g. 120 invocations per day from user center i_1 to s_1 .

$$F = \begin{matrix} & s_1 & s_2 & s_3 \\ \begin{matrix} i_1 \\ i_2 \\ i_3 \end{matrix} & \begin{bmatrix} 120 & 35 & 56 \\ 14 & 67 & 24 \\ 85 & 25 & 74 \end{bmatrix} \end{matrix}$$

The cost matrix $C = [c_{sj}]$, is used to record the cost of deployment of services from candidate locations, which c_{sj} is an integer that indicate the cost of the deployment fee from a candidate location. e.g 130 \$ to deploy s_1 from j_1 .

$$C = \begin{matrix} & j_1 & j_2 & j_3 \\ \begin{matrix} s_1 \\ s_2 \\ s_3 \end{matrix} & \begin{bmatrix} 130 & 80 & 60 \\ 96 & 52 & 86 \\ 37 & 25 & 54 \end{bmatrix} \end{matrix}$$

Consider the following key modeling assumptions:

1. The new WSP decides where to locate his facilities regardless if there is existed functional similar services from other WSPs.
2. This choice is made only consider two factors: total network latency and total cost.
3. We assume a fixed customer allocation policy for WSPs. In practice, Web Services typically offer clients persistent and interactive services, which often span over multiple sessions. Therefore, a dynamic reallocation scheme is not practical as it may disrupt the continuity of the services.

3 Modified NSGA-2 algorithm

NSGA-2 belong to the larger class of evolutionary algorithms (EAs), which generate approximate solutions to optimization and search problems by using techniques inspired by the principles of natural evolution: selection, crossover and mutation.

There are two major differences between standard NSGA-II and modified version. Firstly, in order to avoid repeat evaluation of the Fitness of chromosome which contribute to the majority of computation. After the first generation is initialized, the Pareto front will be stored in the memory. In each generation, when evaluate the chromosome, the chromosome will be check whether it is existed in the memory pool. If so, then the calculation of fitness will be skipped. At the end of each iteration, the Pareto front will be replaced by the latest front. The reason for setting a memory pool is that, after a number of iteration, the population start converging. Most of the population are redundant. Therefore, a memory pool could considerably reduce the repeat computation.

The order major difference is the genetic operators which was discussed in 3.4.

Algorithm 1 modified NSGA-2 algorithm

```

1: Initialize a population of chromosome with random binary value
2: Evaluate with fitness functions
3: Non-dominated sort and assign ranking to each chromosome
4: Evaluate the Crowding distance of each chromosome
5: Initialize the Pareto Front Pool
6: while predefined generation do
7:   Tournament Selection
8:   Crossover
9:   Mutation
10:  Constraint Check and repair operation
11:  for ( doeach chromosome)
12:    if chromosome does not exist in the Pareto front Pool then
13:      Evaluate with fitness functions
14:    end if
15:    Non-dominated sort and assign ranking
16:    Evaluate the Crowding distance
17:  end for
18:  Recombination and Selection
19:  Pareto Front Pool  $\leftarrow$  Current Pareto Front
20: end while

```

3.1 Chromosome Representation

In our problem, the chromosome is the service location matrix A that we mentioned in the previous section.

3.2 Objectives and Fitness Function

The objective functions of this entire problem are following:

- Minimize the total cost of services. n is the number of service, m is the number of candidate location.

$$CostFitness = \sum_{s \in S} \sum_{j \in J} C_{sj} \times A_{sj} \quad (1)$$

- Minimize the network total latency of the services. There are some assumptions we made of network latency. First, we assume the latency = 0 if the service is deployed locally. That is, if the service is deployed in the user center A . Then, the invocation from A to the service will be 0 latency. The second assumption is the latency is symmetrical between cities. e.g. The latency from city A to city B is same with the latency from city B to city A . Thirdly, the total invocation are equally contribute by user centers.
 - For each chromosome, firstly, calculate the number of invocation for each service.

$$Invocation_s = \sum_{i \in I} F_{is}$$

- Calculate the total number of each services.

$$ServiceNo_s = \sum_{s \in S} \sum_{j \in J} A_{sj}$$

- Calculate the number of each service deployed in each location.

$$ServiceNo_{sj} = \sum_{j \in J} A_{sj}$$

- Calculate the average number of each user center contribute to the total number of invocation of each services.

$$AverageInvocation_i = Invocation_s \div ServiceNo_s \div ServiceNo_{sj}$$

- Calculate the latency of each service by multiply latency matrix L by service allocation matrix A (chromosome).

$$LocationLatency_i = \sum_{i \in I} L_{ij} \times A_{sj}$$

- Calculate the total latency of the multiplication of $AverageInvocation_i$ and $LocationLatency_i$. If the service is deployed locally, then the latency equals zero.

Check if deployed locally: $LocationLatency_i = 0$

$$LatencyFitness = \sum_{s \in S} \sum_{i \in I} AverageInvocation_i \times LocationLatency_i \quad (2)$$

The population is initialized based on the problem range and constraints. In our problem, we have two constraints, the cost constraint and minimum number of service constraint. If the initialized chromosome does not satisfy the constraints, then repair operators will attempt to recover from possible constraint violations. The detail of repair operators will be discuss in next section.

3.3 Constraints

The our approach, each chromosome have to satisfy two constraints to be a feasible solution. The first constraint controls the service number which is make sure that each service is deployed in at least one location.

$$\sum_{s \in S} A_{sj} \geq 1 \quad (3)$$

The second constraint is cost constraint which predefined the upper boundary of the total cost. An integer number $CostLimitation$ is defined. This constraint will try to constraints

$$\sum_{s \in S} \sum_{j \in J} C_{sj} \times A_{sj} \leq CostLimitation \quad (4)$$

The constraint is implemented by repair operators which is introduced in the next section.

3.4 Genetic operators

Specific selection, mutation and crossover operators were implemented. It is important to note that mutation and crossover operators can produce solutions that might violate the constraints. Therefore, repair operators are needed to try to maintain feasible solutions. The original NSGA-II use a simulated binary crossover (SBX) Beyer and Deb (2001) and polynomial mutation Raghuvanshi and Kakde (2004) to cope with continuous problem. However, our problem is discretized, therefore we use the regular GA mutation and crossover.

3.4.1 Mutation

The mutation operator works as follows: Initially choose one random location from the chromosome. Then, reverse the bit, e.g. 0 to 1 or 1 to 0.

In this study, a chromosome in a population will be selected based on mutation possibility P_m . As shown in below, a random position will be selected and replaced by a reversed number.



Figure 1

3.4.2 Crossover

The crossover operator in this paper is the single point crossover. The crossover is controlled by crossover probability P_c . The crossover point is created randomly within the length of the chromosome. As in the example below, two parents crossover at a point, then two offspring were generated.

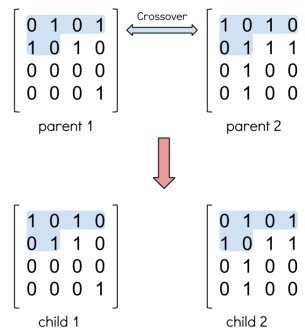


Figure 2

3.4.3 Repair operators

Since mutation and crossover are very likely to generate offspring that violate the constraints, therefore, repair operators are necessary. Normally, each constraint would have a unique repair operator in order to recover different types of violation.

After the process of mutation and crossover, the repair operators are examine each chromosome. If there is violation found in the chromosome, then it will try to repair it. In our problem, we have two repair operators: service number and cost.

Service number repair operator

$$\sum_{s \in S} A_{sj} \geq 1$$

service number constraint

The operator will go through each row of the Location Allocation matrix A, if the number of each service is less than one, then randomly choose one location and reverse the selected bit. Although the randomness does not necessarily provide an optimal solution, the computation time is better than exhaustively compare all solutions.

Cost repair operator

$$\sum_{s \in S} \sum_{j \in J} C_{sj} \times A_{sj} \leq CostLimitation$$

cost constraint

If the cost exceed the predefined limitation. Then the operator would iteratively check if any service has been deployed in more than one location. After found one service has been deployed multiple locations, the operator will randomly select one of them and change to zero. That means, cancel the deployment of that service in this location. After doing that, re-examine the chromosome, if it is still exceed the limitation then repeat this process until there is no redundant services or it satisfies the constraint.

This algorithm will try its best to reduce the cost regardless of other factors. Although the modified chromosome may end up with higher network latency, however the cost constraint has higher priority than network latency.

Worth noting that this algorithm may not provide a strictly correct solution that satisfied the constraint, partially because the randomly chosen of canceling deployment. On the other hand, it may also indicate that the predefined upper boundary of cost limitation is too low.

4 Experiment Design

The purpose of the experiment is comparison between Alternate Location Allocation (ALA) with modified NSGA-II in terms of efficiency and quality of solutions.

The exact algorithm was coded in R using packages: NSGA2R and LpSolve. The program was run on a 3.40GHz desktop computer with 8 GB RAM.

To compare with ALA, four different datasets were used, the problem instance was generated as follows:

1. The potential user center and candidate nodes were randomly selected from existed network latency dataset.
2. The number of potential user center equals the candidate nodes. We defined 4 different size of allocation matrices: 3×3 , 5×5 , 10×10 and 15×15 .
3. Network latency L_{ij} were selected from network latency dataset based on the user centers and candidate nodes that we chose previously.
4. The cost matrix was randomly generate from a normal distribution which mean = 100 and standard deviation = 20

5. The frequency matrix was an integer matrix which randomly generate from a uniform distribution on $[1, 120]$

In each dataset, algorithms run under four different level of cost constraint: Sufficient condition is 100% the expected total cost, good condition (80%), pool condition (40%) and minimum budget condition (0%). The NSGA-II runs 40 times with different random seed ranging from 1 to 40. In efficiency test, We evaluate the average run time for each algorithm.

In order to compare the solution set from NSGA-II and a single solution from ALA, we employ non-dominated Pareto front Xue et al. (2012) and Dominance Relation Ishibuchi et al. (2006). 40 results (from 40 runs) under each cost constraint are presented in the Section 5.2. 40 sets of solutions achieved by each multi-objective algorithm are firstly combined into one union set. In the union set, the non-dominated solutions are presented to compare with the solutions achieved by ALA.

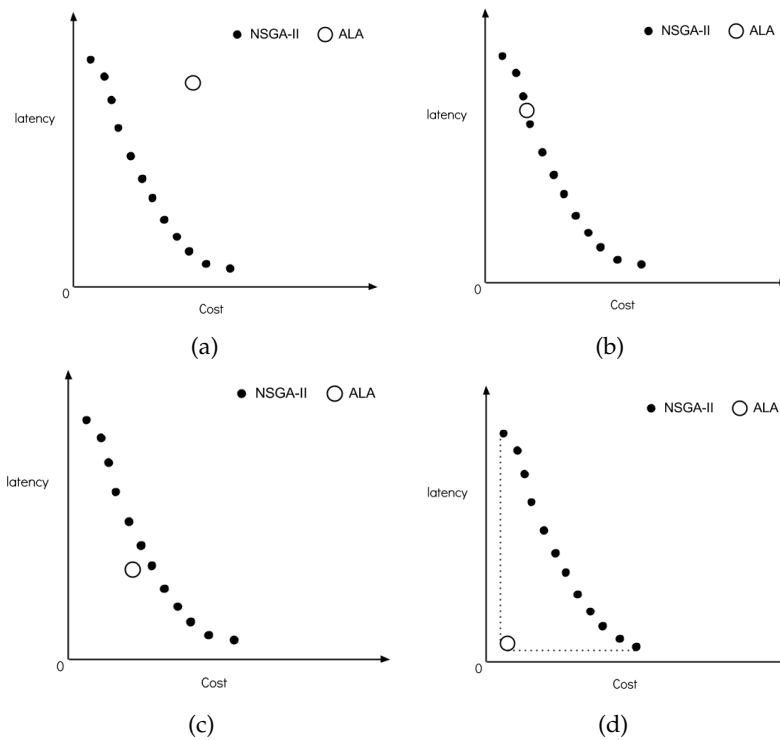


Figure 3: Four cases of the relation between a ALA solution and a NSGA-II solution set of non-dominated solutions

The relation between a single solution and a set of non-dominated solutions from NSGA-II can be classified into four cases as shown in Figure 3. (a), the ALA solution is dominated by some solutions in NSGA-II solution set. In this case, we can say that NSGA-II outperforms ALA since the inclusion of the ALA solution does not improve the quality of the NSGA-II solution set. On the other hand, we can say that ALA outperforms NSGA-II in Figure 3 (d) Since the ALA solution dominates all solutions in the NSGA-II solution set. In Figure 3 (b), the ALA solution is non-dominated solution in the NSGA-II solution set. While the inclusion of the ALA solution somewhat improves

the quality of the NSGA-II solution set, we cannot say that ALA outperforms NSGA-II since the ALA solution dominates no NSGA-II solutions. We may intuitively say that NSGA-II outperforms ALA in the case of Figure 3 (b).

It is difficult to say which is better between ALA and NSGA-II in (c) where ALA solution dominates some NSGA-II solutions. If we use performance measures that are strongly related to the convergence of solutions, ALA maybe evaluated as being better than NSGA-II. On the contrary, if we related strongly to the diversity of solutions, NSGA-II may be evaluated as being better than ALA.

Parameter settings for application of NSGA-II are shown in Table 1. The terminal condition is that the population has evolved 50th generation.

Table 1: NSGA-II Parameters

Population	Generation	Tour Size	Crossover Rate	Mutation Rate
50	50	10	0.8	0.2

4.1 ALA

Much research has been devoted to develop efficient heuristic algorithm to solve location-allocation problem. Two of the predominant greedy construction heuristics are the ADD and DROP Sun and Koehler (2006) methods, which build solutions from scratch. An important greedy improvement heuristic is the Alternate Location-Allocation method which builds on an existing solution. The ADD procedure starts with all facilities closed and locates a facility that gives the greatest savings at each step and iterates through all potential locations until no further savings can be made. In contrast, DROP opens all facilities and removes a location that gives the greatest savings at each step and tries to find optimal or near optimum solutions at the end of the iterative procedure. ALA builds upon ADD and DROP methods. In order to solve multi-objective problem, ALA first uses Integer programming to solve the problem in terms of cost.

$$\begin{aligned}
& \text{Minimize } \sum_{s \in S} \sum_{j \in J} C_{sj} \times A_{sj} \\
& \text{Subject to} \\
& \sum_{s \in S} A_{sj} \geq 1 \\
& \sum_{s \in S} \sum_{j \in J} C_{sj} \times A_{sj} \leq \text{CostLimitation}
\end{aligned} \tag{5}$$

Once a minimum cost solution is found, the solution set will be passed to the ADD procedure to serve as an initial solution for further improvement on Network latency. The ADD procedure as follows:

- Evaluate network latency fitness of the initial solution.
- Starts a facility then evaluate with network latency fitness function. If the network latency fitness decreased as well as the cost fitness does not exceed the cost constraint. The solution is kept.
- If there is no location could be choose, terminate the procedure.

5 Experimental results

5.1 Efficiency comparison

Table 2: Efficiency Test

Table 3: Sufficient condition

Matrix Size	GA time(s)	ALA time(s)	Matrix Size	GA time(s)	ALA time(s)
3×3	1.7	0.002	3×3	1.71	0.003
5×5	2.1	0.007	5×5	2.07	0.01
10×10	3.775	0.11	10×10	3.6	0.11
15×15	6.355	0.64	15×15	6.26	0.6

Table 4: Good condition

Table 5: Poor condition

Table 6: Minimum condition

Matrix Size	GA time(s)	ALA time(s)	Matrix Size	GA time(s)	ALA time(s)
3×3	1.89	0.004	3×3	1.88	0.003
5×5	2.3	0.01	5×5	2.67	0.01
10×10	4.08	0.11	10×10	6.06	0.11
15×15	8.67	0.57	15×15	15.4	0.6

As the experimental results show, in terms of the average run time, ALA is clearly better than NSGA-II. There are two main reasons: firstly, NSGA-II's run time is largely depend on the generation setting. Second, repeat evaluation of chromosome is the decisive factor of time consuming. Although, the modified NSGA-II tries to avoid unnecessary evaluation by employing the idea of memory pool, the improvement is limited.

It is worth nothing that in sufficient condition and good condition, both algorithm's run time remain stable. The run time starts increasing under poor condition and minimum condition with NSGA-II. In particular, the run time for 15×15 matrix under minimum condition is more than twice as under conditions. The reason is that, in the constraint check process, If the children exceed the cost constraint, the repair operator will iteratively close facilities until it satisfied the cost limitation or minimum facility number is achieved. If the cost constraint is set too low, the number of iteration in the repair processes will reach a maximum number. Therefore, the time consuming increases largely.

5.2 Effective comparison

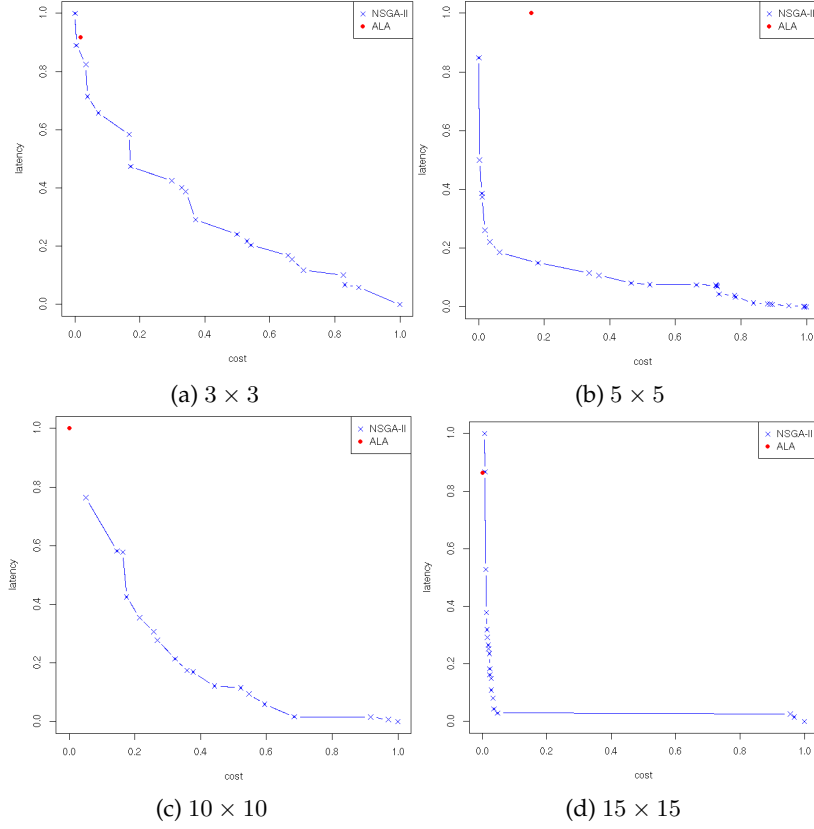


Figure 4

We conducted 16 experiments on 4 datasets under 4 cost constraints: sufficient, good, poor and minimum. The results show similar patterns under same dataset. Therefore, we only pick one figure for each dataset.

Figure 4a and 4b show NSGA-II outperforms ALA. As we discussed in the previous section, the result of ALA is dominated by NSGA-II's solution set. Although, ALA solution for 3×3 matrix is quite close to the Pareto front. 4d shows that solution of ALA is on the Pareto front. It does not improve the quality of the solution set. Therefore, we considered NSGA-II outperforms ALA in this case. 4c is a special case, NSGA-II solution set does not cover the possible solution range. That is, we predefined 50 generations for NSGA-II is not enough for NSGA-II to find a solution set that covers the Pareto front. The solution from ALA is a non-dominated solution in the NSGA-II solution set. While the inclusion of the ALA solution somewhat improves the quality of the NSGA-II solution set, we cannot say that ALA outperforms NSGA-II since the ALA solution dominates no NSGA-II solutions.

As a whole, we can conclude that NSGA-II outperforms ALA in terms of the quality of solutions.

6 Conclusion

In this paper, we presented a modified NSGA-II for web service location allocation problem. Our approach utilizes a memory pool so that greatly reduces the evaluation time while the diversity of the population gradually decreases. We have conducted a full experimental evaluation using the public WS-DREAM dataset to compare our approach to Alternate Location Allocation algorithm. The experimental results shows the modified NSGA-II is effective to provide a near-optima solutions for the web service location-allocation problem.

References

- Aboolian, R., Sun, Y., and Koehler, G. J. (2009). A location allocation problem for a web services provider in a competitive market. *European Journal of Operational Research*, 194(1):64 – 77.
- Beyer, H.-G. and Deb, K. (2001). On self-adaptive features in real-parameter evolutionary algorithms. *Evolutionary Computation, IEEE Transactions on*, 5(3):250–270.
- Caramia, M. (2008). Multi-objective optimization. In *Multi-objective Management in Freight Logistics*, pages 11–36. Springer London.
- Deb, K., Mohan, M., and Mishra, S. (2005). Evaluating the ϵ -domination based multi-objective evolutionary algorithm for a quick computation of pareto-optimal solutions. *Evol. Comput.*, 13(4):501–525.
- Deb, K., Pratap, A., Agarwal, S., and Meyarivan, T. (2002). A fast and elitist multi-objective genetic algorithm: Nsga-ii. *Evolutionary Computation, IEEE Transactions on*, 6(2):182–197.
- Deb, K., Sundar, J., N, U. B. R., and Chaudhuri, S. (2006). Reference point based multi-objective optimization using evolutionary algorithms. In *International Journal of Computational Intelligence Research*, pages 635–642. Springer-Verlag.
- Desai, S., Bahadure, S., Kazi, F., and Singh, N. (2012). Article: Multi-objective constrained optimization using discrete mechanics and nsga-ii approach. *International Journal of Computer Applications*, 57(20):14–20. Full text available.
- Elhossini, A., Areibi, S., and Dony, R. (2010). Strength pareto particle swarm optimization and hybrid ea-pso for multi-objective optimization. *Evol. Comput.*, 18(1):127–156.
- Huang, H., Ma, H., and Zhang, M. (2014). An enhanced genetic algorithm for web service location-allocation. In Decker, H., Lhotsk, L., Link, S., Spies, M., and Wagner, R., editors, *Database and Expert Systems Applications*, volume 8645 of *Lecture Notes in Computer Science*, pages 223–230. Springer International Publishing.
- Huang, V. L., Suganthan, P. N., and Liang, J. J. (2006). Comprehensive learning particle swarm optimizer for solving multiobjective optimization problems: Research articles. *Int. J. Intell. Syst.*, 21(2):209–226.
- Huffaker, B., Fomenkov, M., Plummer, D., Moore, D., and claffy, k. (2002). Distance Metrics in the Internet. In *IEEE International Telecommunications Symposium (ITS)*, pages 200–202, Brazil. IEEE.

- Ishibuchi, H., Nojima, Y., and Doi, T. (2006). Comparison between single-objective and multi-objective genetic algorithms: Performance comparison and performance measures. In *Evolutionary Computation, 2006. CEC 2006. IEEE Congress on*, pages 1143–1150.
- Jamin, S., Jin, C., Kurc, A., Raz, D., and Shavitt, Y. (2001). Constrained mirror placement on the internet. In *INFOCOM 2001. Twentieth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, volume 1, pages 31–40 vol.1.
- Johansson, J. M. (2000). On the impact of network latency on distributed systems design. *Inf. Technol. and Management*, 1(3):183–194.
- Kanagarajan, D., Karthikeyan, R., Palanikumar, K., and Davim, J. (2008). Optimization of electrical discharge machining characteristics of wc/co composites using non-dominated sorting genetic algorithm (nsga-ii). *The International Journal of Advanced Manufacturing Technology*, 36(11-12):1124–1132.
- Raghuwanshi, M. M. and Kakde, O. G. (2004). Survey on multiobjective evolutionary and real coded genetic algorithms. In *Proceedings of the 8th Asia Pacific symposium on intelligent and evolutionary systems*, pages 150–161.
- Ran, S. (2003). A model for web services discovery with qos. *SIGecom Exch.*, 4(1):1–10.
- Sun, Y. and Koehler, G. J. (2006). A location model for a web service intermediary. *Decis. Support Syst.*, 42(1):221–236.
- Vanrompay, Y., Rigole, P., and Berbers, Y. (2008). Genetic algorithm-based optimization of service composition and deployment. In *Proceedings of the 3rd International Workshop on Services Integration in Pervasive Environments, SIPE '08*, pages 13–18, New York, NY, USA. ACM.
- Xue, B., Zhang, M., and Browne, W. N. (2012). Multi-objective particle swarm optimisation (pso) for feature selection. In *Proceedings of the 14th Annual Conference on Genetic and Evolutionary Computation, GECCO '12*, pages 81–88, New York, NY, USA. ACM.
- Zhang, Y., Zheng, Z., and Lyu, M. (2011). Exploring latent features for memory-based qos prediction in cloud computing. In *Reliable Distributed Systems (SRDS), 2011 30th IEEE Symposium on*, pages 1–10.
- Zheng, Z., Zhang, Y., and Lyu, M. (2010). Distributed qos evaluation for real-world web services. In *Web Services (ICWS), 2010 IEEE International Conference on*, pages 83–90.
- Zhou, J. and Niemela, E. (2006). Toward semantic qos aware web services: Issues, related studies and experience. In *Web Intelligence, 2006. WI 2006. IEEE/WIC/ACM International Conference on*, pages 553–557.