

MASARYKOVA UNIVERZITA
FAKULTA INFORMATIKY

Evolutionary optimization of intrusion detection system in wireless sensor networks

DIPLOMA THESIS

Adam Saleh

Brno, spring 2008

Declaration

Hereby I declare, that this paper is my original authorial work, which I have worked out by my own. All sources, references and literature used or excerpted during elaboration of this work are properly cited and listed in complete reference to the due source.

Adam Saleh

Advisor: Andryi Stetsko, Ph.D.

Acknowledgement

I would like to thank my supervisor ...

Abstract

Wireless Sensor networks are an interesting topic for security research because of their possible real world applications, their distributed nature and the constraints on their hardware. For example, setting up intrusion detection system on network of nodes so that it satisfies constraints on accuracy and hardware can be a timeconsuming. We believe that this process could be greatly improved by using genetic algorithms. Purpose of this work is to provide a reasonably reusable framework for optimizations with evolutionary heuristics and calibrate it for optimizing of Intrusion Detection System on Wireless Sensor Network.

Keywords

Wsn,ids,spea2,nsgaii

Contents

1	Introduction	1
2	Evolutionary Optimization of Wireles Sensor Networks	2
2.1	<i>Single objective evolution algorithms</i>	3
2.2	<i>Multiobjective evolution algorithms</i>	4
2.2.1	NsgaII	5
2.2.2	Spea2	5
2.2.3	IBEA	5
3	Our Framework	6
3.1	<i>Evaluation function</i>	6
3.1.1	Bounds	7
3.1.2	Objectives	7
3.1.3	Boinc Assisted Evolution	7
3.2	<i>Calibration</i>	8
3.2.1	Population size and Number of generations	8
3.2.2	Mutation and Crossover	9
4	Experiments	10
4.1	<i>Wireless sensor network</i>	10
4.2	<i>Intrusion Detetection System</i>	10
4.3	<i>Setup</i>	11
4.4	<i>Results</i>	11

1 Introduction

2 Evolutionary Optimization of Wireless Sensor Networks

Optimizing configuration of a wireless sensor network can be a time-consuming task. Nodes of the network are constrained devices, so tradeoffs between battery life, network, memory and processor usage have to be considered. Because it is difficult to predict how will different configuration settings change the performance of the network, it is advisable to measure the impact of changing configuration with WSN simulation software. Unfortunately, simulation of such configuration can take anywhere between several minutes to several hours of processor time, rendering usage of simple exhaustive search on all configurations impractical for most WSN calibration purposes.

According to (Talbi), optimization problems with problem space that is hard to characterize and expensive to search are good fit for heuristic approaches.

[copy figure 1.9 from talbi]

Most of the heuristics used are variants on search through problem space, with different approaches to avoid undesirable local optima.

There are two basic principles that are used when designing a metaheuristic: a) diversification: algorithm should explore the search space

Random search can be considered an extreme application of diversification, where in every round we generate next solution randomly without using any information from previous (good) solutions.

Less extreme example is the family of population based heuristics, where in each iteration we are trying to improve a set of candidates, called a population.

b) intensification: algorithm should exploit information from the best solutions found

Local search can be considered an extreme application of intensification, where next instance to evaluate is always selected deterministically, based on the previous solution.

Local search and its variants belong to a family of single-solution

2. EVOLUTIONARY OPTIMIZATION OF WIRELESS SENSOR NETWORKS

based heuristics, where in each iteration we are basing our search of the problem-space of single solution. (b talbi)

Because of their nature, heuristic approaches are hard to evaluate and reason about, therefore they are often considered to be a method of last resort. (b Talbi) warns against using heuristics lazily on problems where more deterministic approaches (such as using a SAT-solver or linear programming) would be more efficient. Heuristic approaches are usually considered when we need to solve our problem in order of magnitude faster than possible using conventional approaches, while sacrificing guarantees on optimality.

In our case, we are constrained by number of configuration evaluations we are able to perform in a given timeframe.

We are focusing on evolutionary heuristics, because they seem to have reasonable performance with regard of optimizing WSN (b Stetsko 2011). These are modeled natural processes of evolution of species. It is a family of stochastic, population-based heuristics (as opposed to deterministic, single-solution based).

2.1 Single objective evolution algorithms

Basic principles behind evolutionary algorithms can be shown on example of simple evolutionary algorithm.

Evolutionary algorithms are based on the fact, that natural process of evolution can be considered an algorithm solving an optimization problem of adapting a species to its environment.

Template of an evolutionary algorithm. (b Talbi) Generate($P(0)$) ; $t=0$; While not Termination Criterion($P(t)$) Do Evaluate($P(t)$); $P'(t)=$ Selection($P(t)$) ; $P'(t)=$ Reproduction($P'(t)$); Evaluate($P'(t)$) ; $P(t+1) =$ Replace($P(t), P'(t)$) ; $t=t+1$; End While Output Best individual or best population found.

If we continue with this metaphor, we can find these parallels:

Population of individuals is a set of configurations. Their environment is the optimization problem.

We decide surviving individuals based on their evaluation function score.

From surviving individuals we generate next generation, either by mutation or crossover.

2. EVOLUTIONARY OPTIMIZATION OF WIRELES SENSOR NETWORKS

Mutation of individual usually consists of randomly selecting some configuration from its neighbourhood as its offspring. This step provides intensification for evolution heuristic.

Crossover usually consists of creating an offspring by permutation of two individuals. This step provides diversification for evolution heuristic.

In our framework selection and replacement steps are part of the evolutionary algorithm, while crossover and mutation functions are to be calibrated on case by case basis.

2.2 Multiobjective evolution algorithms

Problem that we are trying to optimize is unfortunately ill-suited for single objective evolution. Objectives such as memory consumption and IDS accuracy are orthogonal, therefore it is hard to determine single criteria that would satisfactorily cover both of them. A usual solution is to provide some sort of weighted average. In this case, because output of our algorithm would be only a single solution, we need to know emphasis on different criteria before we run our optimization.

Multi objective algorithms on the other hand are based on idea of Pareto optimality and domination. We say that solution A dominates solution B, if A has no objective "worse" than B and at least one objective "better". We say that A is Pareto-optimal, if there is no other solution that would dominate A. This means that no objective of A can be improved without deterioration of another objective. Set of all non-dominated solutions is then called Pareto front. Set of solutions where no two solutions dominate each other is called Pareto optimal set.

[formal definition of Pareto optimality p 336]

Multiobjective heuristics then try to obtain best approximation of Pareto front. To gauge how good an approximation of Pareto front is, we need to use at least two criteria, first to measure convergence to Pareto optimal front and second to measure diversity in found Pareto set.

[image about approximations]

Calibration of multi objective evolution algorithms is harder than their single criteria counter-parts, mainly because it is hard to reason

2. EVOLUTIONARY OPTIMIZATION OF WIRELES SENSOR NETWORKS

about the convergence of output without having the real Pareto optimal front in the first place. Pareto optimal sets often aren't directly comparable. TODO We will discuss calibration in greater detail in later chapters.

2.2.1 NsgaII

In Non-dominated Sorting Genetic Algorithm, in each generation, population is first separated into ranks based on the ordering of pareto dominance. Then similarity between members of each rank is evaluated. Individuals in population are then sorted, first by rank, second by diversification.

2.2.2 Spea2

In Strength Pareto Evolutionary Algorithm

2.2.3 IBEA

3 Our Framework

(b Stetsko 2011) showed, that using simple evolutionary algorithms is a viable and efficient way of optimizing Wireless Sensor networks. Unfortunately, previous experimental setup was too tightly coupled with optimization framework, rendering the code of previous experiments hard to reuse. Therefore we have created a micro-framework based on python and ParadisEO.

ParadisEO is a highly configurable framework for metaheuristics written in C++. Our micro-framework simplifies its usage by allowing to specify evaluation, mutation and crossover functions in python. Emphasis is on multiobjective algorithms, ability to evaluate population in parallel and having good facilities for experiment analysis.

We chose python, for these reasons:

- 1) it is an efficient glue language. Most often you will want to optimize already written application with minimal changes. Python is well suited for running external applications, specifying their inputs and parsing their outputs.

- 2) it has good facilities for statistical analysis. That allows us to include analysis of every optimization run into the executable itself. Having automatically generated experimental log has proven invaluable especially in calibration of algorithms for our specific problem.

- 3) it is easily integrated with C++. We strive to make the integration with ParadisEO simple and extensible.

3.1 Evaluation function

To simplify configuration, inputs and outputs of evaluation function are always a list of floating point numbers. We believe this is general enough, and that most WSN optimisation problems can be fitted to this constraint.

3.1.1 Bounds

For inputs you specify their bounds, which are then used to initialize the population and to supply constraints for mutation and crossover parameters

[Example]

3.1.2 Objectives

Because we focus on multiobjective evolution, you need to specify number of objectives and which ones we want to minimize (maximize).

[example]

It is advisable to have three objectives at maximum, because number of solutions on the Pareto optimal dramatically increase with new objectives. At minimum, $(n+1)$ objective problem will contain all the solutions of an n -objective problem. [CITATION? 309 talbi]

Then there is the problem of visualizing and evaluating more than 3-dimensional data.

3.1.3 Boinc Assisted Evolution

Because evolution is population-based heuristic, it is well suited for parallel evaluation of its individuals. ParadisEO already has two methods included, either by use of MPI protocol [CITATION], or by using SMP on multi-core machine. Unfortunately, SMP module of ParadisEO is not yet stable on all platforms, and MPI has specific demands on infrastructure and is usually hard to retro-fit on already written program.

Boinc on the other hand can simply utilize machines already present to form a simple computational grid. With its simple architecture consisting of Management server and several worker nodes it can be deployed in most settings. Its simple architecture prohibits cooperation between worker nodes, but that doesn't concern us while evaluating individuals.

We have decided to give the evolution algorithms ability to use a simple scheduler for creating Boinc work units, pluggable from python.

In our micro-framework it consists of three functions, `schedule_work_unit`,

`wait_for_completion` and `gather_result`.

Even though the evaluations themselves will be done in parallel, we wanted to avoid any threaded code in our framework itself. Because we have written it with research in mind, we understand that accessing auxiliary data is important to evaluate experiments. In threaded environment, this could lead to deadlocks.

3.2 Calibration

We can think of calibration on an evolutionary algorithm as a meta-optimization problem with two objectives,

1) maximizing the fitness of result 2) minimizing the number of evaluations

With single-objective algorithms, optimizing for these two objectives is easy, because both of them are usually represented by numbers. With multi-objective problems, results are Pareto sets, which are often not directly comparable.

Quality of a multiobjective result can be judged based on several criteria:

1) number of results in Pareto set. 2) distance from true Pareto front 3) regularity in distance between results

Because the main reason we use multi objective algorithms is to avoid the need to do the expensive computation of true pareto front, we could either calibrate on a smaller sample of the problem

Underlying assumption behind sampling is, that the Pareto front of the solutions of the problems subset will be the subset of true pareto-front.

Or we could use different convergence criterion. [TODO]

This works well only if we can assume that the solution-space is homogenous, can accurately gauge its boundaries.

3.2.1 Population size and Number of generations

These two parameters set the boundaries on minimal and maximal number of evaluations. Lets mark the size of population N and number of generations G Because the whole initial population needs to be evaluated, number of evaluations will at least the size of populations.

Because at each new generation, at most N new individuals are generated, therefore in one optimization run, at most $N \times G$ evaluations.

3.2.2 Mutation and Crossover

Probability of mutation, or crossover directly influence number of generated individuals in each generations.

4 Experiments

To test our framework we have decided to optimize a simple intrusion detection system on a model wireless sensor network.

4.1 Wireless sensor network

Implementation of WSN was reused from previous research (Stetsko 2012), with just a slight modifications to decouple it from previous optimisation framework. Nodes are simulated with MiXiM framework, with test application on each of them, that sends a packet containing arbitrary information through the network to a base station. Network has a static routing tree based on [TODO] algorithm.

We have marked some of the nodes as "droppers", so that instead of forwarding all of the packets through the

4.2 Intrusion Detection System

Our nodes implement a simple IDS capable of detecting nodes that seem to be intentionally dropping packets. If we look at the wireless network stack of our node, IDS is part of the medium access control sublayer. There it can eavesdrop on neighbouring nodes and check whether they behave accordingly.

For each neighbour it updates a) number of packets received b) number of packets forwarded.

We wanted to use IDS that is simple, but highly configurable. We have these four parameters to optimize: 1) number of monitored nodes 2) size of buffer for eavesdropped packets 3) threshold for minimal number of received packets for a node to be considered malicious 4) threshold for ratio between forwarded and received packets for a node to be considered to be malicious

At the end of simulation, based on a preset thresholds and ratio between forwarded and received packets it decides whether node is malicious or benign. More specifically, the set of malicious nodes is:

And the set of benign nodes is:

You can see, that both sets are disjoint, and their union is the set

of neighbours N

4.3 Setup

4.4 Results