

Master 1 MoSIG Research Project Report

Learning Job Runtimes in Homogenous HPC Systems

Valentin Reis

Supervised by: Denis Trystram & Eric Gaussier.

I understand what plagiarism entails and I declare that this report is my own, original work.
Name, date and signature:

Abstract

Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua. At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.

1 Introduction

High Performance Computing (HPC) systems are complex machinery at the frontier between research in scheduling and systems engineering. We outline two of the main difficulties that resource management software in this field have to face.

First, the ephemeral nature, and broad range of existing architectures of such systems make the development and application of theoretical results difficult. New schemes for distributing resources are more present than ever. Many systems have complex network and memory/hard drive sharing topologies. Moreover, the topology of some HPC system can change on a hourly to monthly basis, as the supporting hardware of distributed systems can be reconfigured or extended continually. Finding scheduling and resource management strategies which can deal with complex systems and adapt to their evolutions poses a challenge.

In addition, the data (i.e. the characteristics of the tasks to be executed on the system) these systems have to work with presents many peculiarities. The nature of the information which users of the system provide is very often loose (e.g. with only upper and/or lower bounds on numerical quantities provided). For instance, and this will be the focus of this paper, the run time of a given job on a specific system is seldom known in advance, but many cluster management software ask the users for an upper bound on this quantity.

As a consequence of these difficulties most free, open-source and commercial resource management software use simple heuristics that can provide bounds on their performance and/or guarantee a few functional properties. An example of such a heuristic is the First Come First Serve (FCFS) policy to schedule parallel jobs on a homogenous cluster of

machines, which starts jobs as soon as possible, in the exact order they were submitted. Among other properties (such as robustness to lack of information about the amount of time jobs will run on the system), this strategy guarantees the avoidance of starvation.

1.1 Research Direction

The general direction we are headed in with this research project is to deal with the input data of the resource management systems. Accomodating for this data seems separable enough from the actual scheduling problems for work towards this objective to be rewarding. No innovative ways to query the data from the users will be studied : we will mainly rely on existing logs from HPC systems. Instead, we seek to apply Machine Learning techniques in order to reduce uncertainty of, and extract information and/or structure from, the input data of the HPC systems. We will be working with the problem of presenting input data in the most valuable way possible to a scheduling algorithm. How to use this data to the fullest will not be discussed. When assessing the relevance of the specific information we choose to produce from the job characteristics, references from the scheduling literature and existing systems will provide ground to stand on.

1.2 Job runtimes

Most HPC resource management software (including the SLURM, OpenPBS, OpenLava and OAR software) do ask information about jobs to users, such as topological requests in terms of processing units and memory, the name of the executable, miscellaneous functional requirements and, last but not least, the expected run time of the job with respect to its hardware requirements. This user-provided estimate of the run time of a job on a specific system will be referred as **reqtime** in the rest of the paper. Most of these software use the **reqtime** of a job as an upper bound on its run time, and kill it should **reqtime** be violated. As a consequence, users overestimate this value, should they choose or be forced to provide it. The following section will present a statistical analysis pertaining to this relationship.

The true run time of a job with respect to a given affected topology is of great interest, as the scheduling policies are highly dependent on this information to provide good solutions. We will refer to this quantity as the **runtime** of a job. It must be clear that in the context of topological hetero-

geneity, the **runtime** of a job is only defined with respect to a specific processing environment to which it might be affected. This can include, and is not limited to, the network topology of the processing units, the availability of shared memory, message passing costs, and the operating system supporting the computations.

1.3 Problem Statement

In this paper, the broad question of refining the data is reduced to a single variable, the **runtime**. The problem statement we are dealing with is the following.

Given a specific homogenous HPC Cluster with negligible communication costs, how to best predict the value of the **runtime** of a job?

The choice of dealing with homogenous distributed machines without communication costs in a first approach has the interesting property of separating the data treatment from the scheduling and interaction with the system. In this case, the **runtime** becomes an intrinsic attribute of a job. On the contrary to the input data which is always subject to peculiarities of the various systems and software, the **runtime** is always present as a simple field in workload logs of HPC systems, such as those available at the workload archive [ref to feitelson].

This problem statement implies a latent question : How to communicate the prediction to the rest of the system (e.g. single-value, probability density, confidence factor)? Alternatives will be discussed but ultimately, the focus will be on single-valued predictions. As mentioned previously, our approach is to use machine learning techniques. We will be inferring runtime from the the job characteristics by learning from system logs, and since this is a value in $[0, +\infty]$, this is a standard supervised learning problem, namely regression.

We will be careful to only learn our models on logs from homogenous systems, or homogenous subsets of systems which are sizeable enough to learn from. As for downplaying the impact of communication costs, we will further restrict our work to machines which do not possess overly complex topology or distribute computing nodes across more than a handful of routers. In essence, we are targeting large Beowulf clusters, supercomputers, GPU farms and mainframe clusters.

2 Motivation

2.1 Importance of runtime

Once again, we emphasise the role of the **runtime** in scheduling tasks. Classical results from scheduling theory use the "clairvoyant" model[ref], where all **runtimes** of jobs with respect to any affectation on the system being modeled are known. Informally, if we want to make use of this extensive theoretical body, we must at least reduce uncertainty in this information. More recent models [ref Trystram FCFS] do [...]. Another approach [ref prob.backfilling] uses predicted probability distributions on the runtimes of jobs in order to take scheduling decisions. In all cases, reducing uncertainty in **runtime** is critical to the success of the approach used to schedule jobs.

As mentioned before, existing solutions use much simpler heuristics. We will now focus on a particular system, and outline the reasons why, in the absence of a predictive technique applied the input data, the FCFS heuristic is applied.

2.2 reptime vs runtime on a real system

The following study is conducted on a log[ref log] containing 20 months worth of data from the Curie[ref curie] supercomputer operated by the French government-funded research organization CEA (Commissariat à l'Énergie Atomique). The system has several highly homogenous partitions, to which jobs are allocated according to user preference. Figure [ref figure] shows the marginal distributions of **reptimes** and **run-times** on this system.

TODO: show the curie log.

3 State of the art in runtime prediction

3.1 Nature of the prediction

TODO: -what to predict: value?, confidence interval?, distribution? which algorithms can we use those?

3.2 Predicting a value

TODO: -give the references and explain the historical methods, gibbons historical scheduler, the tsafir et al paper with mean of two last runtimes values userwise..

3.3 Predicting a distribution

TODO: -give the references and explain the probabilistic backfilling thing..

4 Our approach

WHY ML? well, there is simple little patterns everywhere

4.1 Random Forests

TODO: -explain our approach, why it could lead to better results (external info+signal locality(ref hmm thesis for locality..))

4.2 Explainability

TODO: -explain one advantage of random forests: discussion about the trees after learning..

5 Preliminary Results

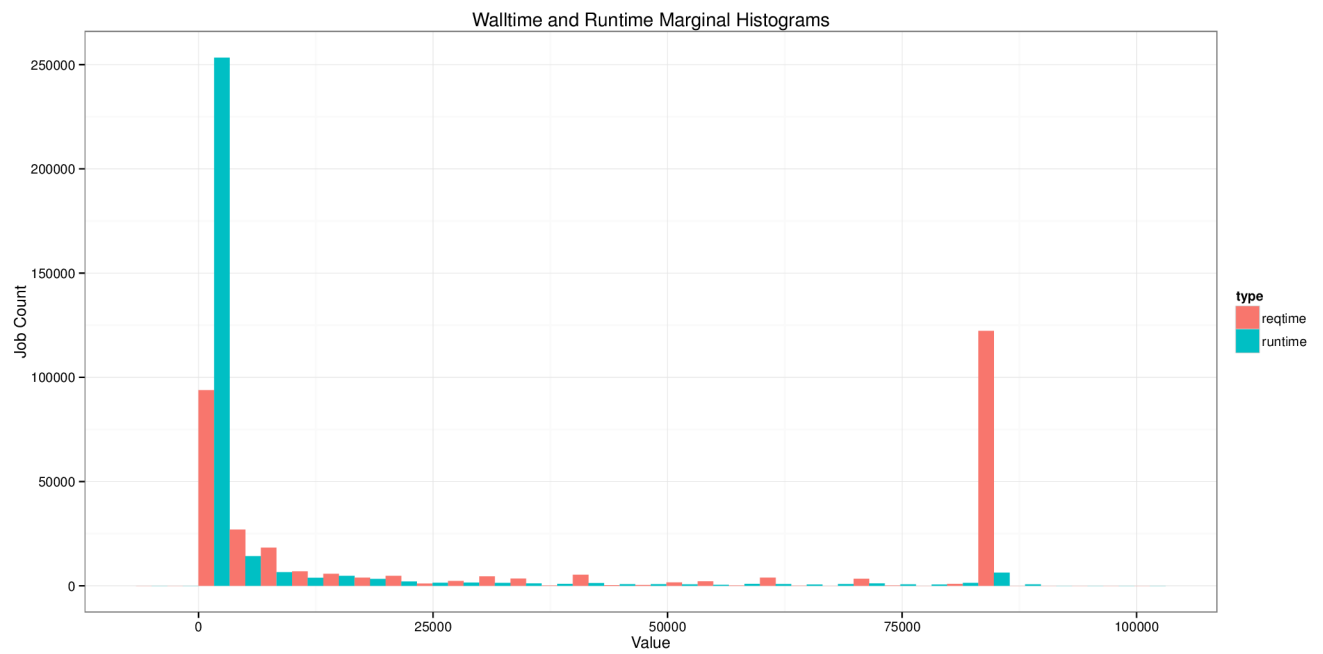
TODO:results..

6 Conclusions

TODO:conclure..

7 Acknowledgements

TODO:remercier..



/3/3

Figure 1: Marginal Histograms of **reqtime** and **runtime** of jobs on the CEA Curie log.