

Noah L. Schrick  
Professor: Dr. Abdulrashid  
Final Project Proposal

## **Topic: Data Mining Techniques on Algorithm Analysis for Compliance Graph Generation**

### **Introduction**

Compliance graphs are generated graphs (or networks) that represent systems' compliance or regulation standings at present, with expected changes, or both. These graphs are generated as directed acyclic graphs (DAGs), and can be used to identify possible correction or mitigation schemes for environments necessitating compliance to mandates or regulations.

Compliance graphs are an appealing approach since they are often designed to be exhaustive: all system properties are represented at its initial state, all violation options are fully enumerated, all permutations are examined, and all changes to a system are encoded into their own independent states, where these states are then individually analyzed through the process. Despite their advantages, compliance graphs do suffer from their exhaustiveness as well. Even very small networks with only 10 hosts and 5 vulnerabilities yield graphs with 10 million edges. When scaling compliance graphs to analyze the modern, interconnected state of large networks comprising of a multitude of hosts, and utilizing the entries located in the National Vulnerability Database and any custom vulnerability or regulation testing, graph generation quickly becomes infeasible. Similar difficulties arise in related fields, where social networks, bioinformatics, and neural network representations result in graphs with millions of states. This state space explosion is a natural by-product of the graph generation process, and removing or avoiding it entirely undermines the overall goal of compliance graphs.

Related research works extend the generation process to function on distributed computing environments to take advantage of the increased computing power using message-passing. The approach makes use of a task parallelism approach for the generation process, and examines the effect of various parameter tuning and algorithm design choices and their effects on speedup and efficiency.

### **Data**

The algorithmic approach is to break the generation process into six main tasks. Dependent variables of the data are time per task and overall runtime. To measure the speedup and efficiency, data is collected for each task based on parameter tuning, as well as data for the overall process. Reducing the dependent variables, it can be shown that there is one dependent variable (overall runtime) that is comprised of 6 sub-dependent variables (time per task).

Across the six main tasks and throughout the generation process, a total of 4 additional independent variables are used. These independent variables are: the number of compute nodes used, the number of exploits, the number of applicable exploits, and the memory load before database operations are performed. Each variable is changed individually across a given range to capture a full-spectrum image of their effects on time.

The number of nodes ranges from 2 through 12, incrementing by 1.

The number of exploits ranges from 6 to 49152, incrementing by a factor of x2.

The number of applicable exploits ranges from 0% to 100%, incrementing by 25%.

The database load ranges from 0% to 100%, incrementing by 25%.

The dataset therefore has 2,464 entries, with each entry spread across 11 columns, leading to a total of 27,104 individual pieces of data. The data will be collected in an automated fashion. The experiment will be deployed on the University's supercomputer cluster, with custom scripting to automatically run each subsequent test, and collect the data into a raw, unfiltered CSV.

## **Methodology**

### *Analysis*

The methodology will work to answer the following questions:

1. How does each parameter contribute to the overall runtime?
2. Which parameters have the largest effect on overall runtime?
3. Can an optimal set or subset of parameters with corresponding values be identified?
4. What is the best way to visualize and represent the data and findings?

The methodology will utilize multivariate analysis, covariance, sums of squares, regression, heatmaps, and will leave room to explore methods such as PCA or LASSO.

### *Results*

Ideal results will include definite answers to analysis questions 1, 2, and 3 through figures and various discussion. Results for analysis question 4 will include figures along with a discussion on whether the figures successfully capture all findings.

### *Appropriateness of Approach*

Since data is to be collected only for the effect of the parameters on overall runtime, additional care is needed to ensure that noise is filtered as best as possible. Timings for the algorithm need to be carefully inserted so as to not bias or detract from the goal of the research. Substantial discussion will need to be included to justify the regression approaches.

### *Conclusions*

Ideal conclusions will have definite answers to analysis questions 1 through 3, with substantive figures and discussion for analysis question 4.

## **Limitations of Analysis/Data**

A discussion and justification will be required for both regression and PCA and/or LASSO. PCA and LASSO will both work toward dimensionality reduction, but achieves the goal in separate means. A discussion of limitation will be required to best answer which approach best fits the data, and shortcomings of each approach as it relates to the data.

## **Future Works**

Based on the results, future work could include discussion of future algorithm improvements, suggestions for corrections or alternate approaches on algorithm design or analysis approaches, or highlight novel or unseen effects of parameter tuning as they relate to speedup and efficiency of compliance graph generation.