

並列分散システムにおける異常動作の原因特定のためのログ解析*

佐伯 勇樹†

田浦 健次朗‡

近山 隆§

東京大学工学部電子情報工学科† 東京大学大学院情報理工学系研究科‡ 東京大学大学院新領域創成科学研究科§

1 背景と目的

近年、大規模な計算を可能とする並列分散計算が注目を浴びているが、プログラムのバグや設定ミス、故障などたびたび異常が発生する。これらの異常は物理的に離れたノードや時間的に離れたイベントに影響を与えることもあり、その原因特定は既存のシステムでも完全に行うのは困難であるため、そのような異常の原因を特定するための手段を用意する必要がある。そこで、情報爆発時代に向けた IT 基盤技術研究のための研究プラットフォーム“InTrigger”をテストベッドとし、機械学習を用いて、各ノードの CPU 使用率や通信量などのログからノードの状態分類を行い、そのノードの接続情報からどのプロセスによる異常かを特定する手法の提案と実装、そしてその評価の報告をする。

2 関連研究

既存のネットワークモニタリングシステムとして ganglia [4] が上げられる。このツールではプロセス一つ一つのデータに関しては情報を得ることができず、また統計量だけからでは異常原因特定につながりにくいことが多い。またネットワーク上の Web サーバ数台による分散環境の状態を分析する方法として、Cohen らによって Tree-Augmented Naive Bayes (TAN) を用いた手法が提案されている。Cohen らはサーバの応答時間を正常と異常の分類基準とし、TAN [3] を用い、サーバの CPU 時間、ネットワーク量などの統計量から、状態推定を行った [1]。そして調査したいシステムの各統計量を TAN で 2 値に特徴づけ、クラスタリングする実験を行い [2]、異常状態の特徴を見つけることに成功した。問題としては数百といった大規模環境ではベクトルの次元数が大きくなりすぎることが挙げられる。

3 本研究の手法

本研究は、各マシンから得た統計量やプロセス情報でのクラスタリングを大規模環境で行う方法を提案し、その結果からマシンの異常状態の種類を調べ、その原因を状態分類とそのプロセス情報から推測する。

3.1 実験環境と解析を行うデータ

今回評価する並列分散システムは、前述の InTrigger を用いた。解析を行うデータとして、マシンの状態を表す CPU 使用率、ディスクアクセス回数といった各種統

計量などを /proc, vmstat, ps, netstat, iptraf から取得しログに記録する。

3.2 手法の概要

状態のクラスタリングは、ある一時点の 1 分間の統計量のベクトルをクラスタリングすることで行うが、各々のマシンは動いているプロセスにより統計量の変動に特徴があり、正常でも変則的な値をとったり、異常でも一見正常な値に見えることがある。そこで、Cohen らの 2 値変換の手法を大規模な環境向けに修正を加えて用い、異常原因となっている統計量を明確にする。まず、2 値変換にはナイーブベイズ、TAN の 2 種類の学習構造を事前に蓄えたログから学習し、適したものを使う。

ここで用いるナイーブベイズとは、各々の変数がクラスに依存しており、ある時点での統計量ベクトル x が m であったときに、その状態クラスを最も確率 $P(s|x=m)$ の高いクラス s に分類する手法である。その確率は次のようになる。

$$P(s|x=m) \propto \prod_{i=1}^n P(x_i = m_i | s) P(s)$$

ここで確率分布 $P(x_i = m_i | s)$ は正規分布を仮定し、学習データの平均、分散から求める。

TAN [3] とはナイーブベイズの学習構造において、各々の変数はクラス以外に最大 1 つの異なる変数に依存しているとする学習構造である。ただし依存関係に閉ループはないものとする。これは、各統計量間の依存関係の強さを表す相互情報量をエッジ、各統計量成分をノードとするグラフ G の最大全域木を求める問題に帰着できる。

手順は次のようになる。

まずは学習用ログデータを、ssh 応答が速いか遅いか、ping 応答が速いか遅いか、loadaverage が変に大きいか正常な値であるかなど、閾値によって単純に正常クラスと異常クラスの 2 つに分ける。このクラス分類をナイーブベイズ、TAN で学習する。

次いで調べたいログデータを 2 値変換する。統計量の 2 値変換に使う学習構造の選び方は、分類したい時点 t の近辺 w の正常、異常ラベルをよりはっきり推定できる学習構造を Brier Score という基準によって判定して選び、2 値変換は選んだ学習構造での単純な確率比較によって、正常時にとる確率が高いなら +1、そうでない場合 -1 にする。

* Log Analysis for Finding Trouble Sources in Distributed Systems.

†Yuki Saeki, Department of Information and Communication Engineering, Faculty of Engineering, The University of Tokyo

‡Kenjiro Taura, Graduate School of Information Science and Technology, The University of Tokyo

§Takashi Chikayama, Graduate School of Frontier Sciences, The University of Tokyo

さらに調べたいシステムの状態群を各要素が 2 値のベクトルで表したら、これをクラスタリングする。ここでは kmeans 法を採用した。ある程度クラスタリングされたデータが集まった後は、各クラスタの中心を求めておき、以降の分類はどの中心に近いか、もしくはどの中心ともかけ離れている新たな異常であるかで状態分類することができる。状態分類を行う環境を 50~100 台の並列分散環境に適応させるため、ここでは同一時点のそれぞれのマシンの状態を別々に集める。時刻 t の状態データは分散環境のマシンの台数分集まることになり、それを調べたい区間で集め、まとめて kmeans クラスタリングを行う。また、NFS サーバと計算ノードは別々に kmeans を行うとする。

最後に状態分類をした後は、netstat コマンドなどのネットワーク接続情報、ps コマンドの出力などのプロセス情報からその状態とかかわりの深いアクティブであるプロセスを抽出し、原因を判定する。具体的には、プロセスグループごとにプロセスをまとめ、同一プロセスグループ ID を持つプロセスが極端に多い(50~100 以上)グループ、CPU 使用率やメモリ領域に変動のあるグループ、STAT が R、D になっているプロセス数の多いグループ、ネットワーク通信量の多いグループなどである。

4 評価

今回の評価の一例として、84 の計算ノードと 1 つの NFS サーバからなる hongo クラスタ環境で起こった異常を取り上げる。

1 月 6 日から 8 日の間、あるユーザのプロセスの異常により、NFS サーバに大きな負荷がかかった。この区間でのデータを収集し、クラスタリングを行った。

クラスタリングに用いた統計量は、メモリの空き容量、割り込み回数、コンテキストスイッチ回数、CPU ユーザ時間、CPU システム時間、IO 待ち時間、ディスク読み込み回数、ディスクアクセス回数、ネットワークデータ送信量、ネットワークデータ受信量の 10 項目で、正常と異常の大まかな分類基準として loadaverage を閾値 5.0 で用いた。

まず学習データとして計算ノード一つと NFS サーバの正常と異常の混じった 1 日分のデータを学習させ、続いて、この学習データを用い、すべてのノードの 3 日間のデータを 2 値変換し、ノード群と NFS サーバとで別々に kmeans 法を計算ノードでは $k=7$ で、NFS サーバでは $k=10$ で実行した。すると、次のことがわかった。

- 通常時は計算ノードが状態 2 か 6、NFS サーバが状態 4 になった。
- 6 日の 15:30 近辺で 1~69 番の計算ノードが状態 3 に、0 番の計算ノードが状態 0 と 3 をふらつく状態になった。
- 6 日の 19:00 近辺で 1~69 番ノードは状態 2 に戻るが、0 番ノードのみ状態 4, 5 をふらつく状態へと遷移した。
- 7 日から 8 日で、0 番ノードは状態 4, 5 をふらつく状態と状態 1, 4, 5 をふらつく状態を繰り返した。

- 7 日から 8 日で、NFS サーバは状態 0, 2, 3, 7, 8 をふらついていたが、時間帯によって状態 2 が多い場合と状態 8 が多い場合が存在した。それぞれ 0 番ノードが状態 4, 5 をふらついていたときと状態 1, 4, 5 をふらついていたときにに対応していた。

プロセス情報を調べると、6 日 15:30~19:00 に 0~69 番ノードで同一名のプロセスが立ち上がりつづいていたが、うち 0 番ノードのみそれ以降もそのプロセスが残り続け、30 以上に fork し IO 待ち状態になっていた。また状態の遷移に伴い IO 待ちの数が変動していた。さらに、NFS サーバでは NFS デーモンがそれに対応して大きく負荷をかけているのがわかった。

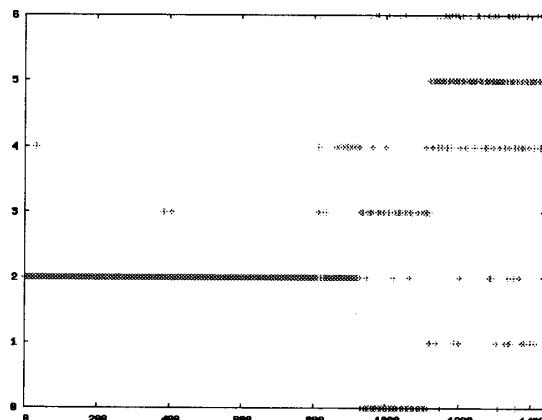


図. 0 番ノードでの 1 月 6 日のクラスタリングの結果

5 まとめと今後の課題

各ノードごとの状態を別々に集めてクラスタリングすることにより、マシンの状態遷移と原因プロセスを追いかけることを可能にできたが、状態のふらつきが多いために、完全な自動判断にはまだ改良の必要がある。分類に用いる統計量の選択、原因プロセスの学習、十分な学習量をえることが今後の課題である。

参考文献

- [1] I. Cohen, M. Goldszmidt, T. Kelly, J. Symons, J. S. Chase. Correlating instrumentation data to system states : A building bloc for automated diagnosis and control. In Proc. 6th USENIX OSDI, San Francisco, CA, Dec. 2004.
- [2] I. Cohen, S. Zhang, M. Goldszmidt, J. Symons, T. Kelly, A. Fox. Capturing, Indexing, Clustering, and Retrieving System History. In SOSP'05.
- [3] N. Friedman, D. Geiger, and M. Goldszmidt. Bayesian network classifiers
In Machine Learning, vol. 29, pp. 131 163, 1997.
- [4] Ganglia Monitoring System.
<http://ganglia.sourceforge.net/>