

Proactive Fault Management Guest Lecture in "Dependable Systems" by Dr. Peter Tröger

Dr. Felix Salfner

19.6.2012 felix.salfner@sap.com



Contents

Introduction

- Variable Selection
- Online Failure Prediction Overview
- Four Online Failure Prediction Techniques
- Assessing Failure Predictors
- Taking Action
- Summary



Our Credo

"Ordinary mortals know what's happening now, the gods know what the future holds because they alone are totally enlightened.

Wise men are aware of future things just about to happen"

C. P. Cavafy, (Greek poet, 1863-1933) "But the Wise Perceive Things about to Happen," a poem based on lines by Philostratos



Motivation

- Ever-increasing systems complexity
- Ever-growing number of attacks and threats, novice users and third-party or open-source software, COTS
- Growing connectivity and interoperability
- Dynamicity (frequent configurations, reconfigurations, updates, upgrades and patches, ad hoc extensions), and
- Natural and man-made disasters



A Different Mindset

- Faults, errors and failures are common events so let us treat them as part of the system behavior and learn how to cope with them
- Attractive panacea:

(self) Proactive Fault Management (PFM)

Proactive Fault Management (PFM)

PFM is an umbrella term for techniques such as monitoring, diagnosis, prediction, recovery and preventive maintenance concerned with proactive handling of errors and failures: if the system knows about a critical situation in advance, it can try to apply countermeasures in order to prevent the occurrence of a failure, or it can prepare repair mechanisms for the upcoming failure in order to reduce time-to-repair.





Comparison to Classical Reliability Theory

- Classical reliability theory is typically useful for long term or average behavior predictions and comparative analysis
- Classical reliability theory may help but is not very good for short term prediction due to dynamics, mobility, systems/networks complexity, changing execution environments, upgrades, online repair, etc.



Contents

- Introduction
- Variable Selection
- Online Failure Prediction Overview
- Four Online Failure Prediction Techniques
- Assessing Failure Predictors
- Taking Action
- Summary



Variable Selection

- What are the right variables to use for modeling?
- There are up to 4200 variables (v) and up to hundreds of fault classes (f) per node
- For n nodes: m = v x f x n variables, the number of combinations c equals:

$$\mathbf{c} = \sum_{r=1}^{m} \binom{m}{r}$$

• Combinatorial explosion!



- Selection by experts
- Filter (e.g., mutual information criterion)
- Wrapper (making use of modeling procedure specifics)
 - feed forward selection, finding independent variables
 - backward elimination

- probabilistic (only variables showing correlation and certain distribution)

 Forward Addition - a method of selecting random variables for inclusion in the regression model by starting with no variables and then gradually adding those that contribute most to prediction



Variable Selection



PFM - 19.6.2012 - 11



Contents

- Introduction
- Variable Selection
- Online Failure Prediction Taxonomy
- Online Failure Prediction Techniques
- Assessing Failure Predictors
- Taking Action
- Summary

Faults, Errors, Failures again ...



PFM - 19.6.2012 - 13

Four Ways of Detecting Faults

(1) The system can be *audited* in order to actively search for *faults*, e.g., by testing on checksums of data structures, etc.

- (2) System parameters such as memory usage, number of processes, workload, etc., can be *monitored* in order to identify side-effects of the faults. These side-effects are called *symptoms*. For example, the side-effect of a memory leak is that the amount of free memory decreases over time.
- (3) If a fault is activated and *detected* (observed), it turns into an *error*.
- (4) If the fault is not detected by fault detection mechanisms, it might directly turn into a *failure* which can be observed from outside the system or component.



Taxonomy





Online Failure Prediction -Definition

- The goal of online failure prediction is to *identify failure-prone situations*, i.e. situations that will probably evolve into a failure. The evaluation is *based on runtime monitoring data*.
- The output of online failure prediction can either be
 - a decision that a failure is imminent or not, or
 - some continuous measure evaluating how failure-prone the current situation is

Two Types of Input Data

- There are two types of system measurements
 - periodic, numerical
 - event-based, categorical
- Examples for periodic data
 - system- / CPU load
 - memory usage
- Examples for event-based data:
 - interrupts
 - threshold violations
 - error events





Contents

- Introduction
- Variable Selection
- Online Failure Prediction Overview
- Four Online Failure Prediction Techniques
- Assessing Failure Predictors
- Taking Action
- Summary



- 1. Universal Basis Functions (UBF)
- 2. Hidden Semi-Markov Model (HSMM)
- 3. Dispersion-Frame Technique (DFT)
- 4. Eventset method

, Universal Basis Functions (UBF)



- Tailored to periodic measurements
- Function approximation approach: Express target value as function of input variables
- Examples for target values:
 - Availability
 - Memory consumption



UBF Background

Starting from radial basis functions
Linear combination of kernel functions G_i

$$f(\mathbf{x}) = \sum_{i=1}^{n} \alpha_i * G_{\lambda_i} (\|\mathbf{x} - \mathbf{c}_i\|) \qquad \text{Replace domain} \\ \mathbf{f} = \mathbf{c}_i \| \mathbf{x} - \mathbf{c}_i \|$$

Replace fixed Gaussian by flexible domain specific kernel

Determine parameters by minimizing, e.g., mean square error

$$\min H[f] = \sum_{i=1}^{N} \left(f(x_i) - y_i \right)^2$$

PFM - 19.6.2012 - 21

Effect of ω (UBF slider)

- Linear combination of nonlinear kernel functions
- Examples:
 - Gaussian
 - Sigmoide functions, ...
- RBF is special case
- Improve efficacy by introducing data specific flexible kernels
- Universal approximator
- Large number of kernels to cover heavy tailed distributions



PFM - 19.6.2012 - 22

Dispersion Frame Technique



• Classic technique for error-log analysis

- Evaluates the time of error occurrence
- Applies a set of heuristic rules evaluating the number of errors within successive dispersion frames

[Lin, Siewiorek 1990] PFM - 19.6.2012 - 23



Event-set Method

- Approach inspired by data-mining
- Focus on type of events
- Based on sets of events
 - Each set contains decisive events that occur prior to a target event
 - Events correspond to errors in our taxonomy
 - Target events correspond to failures
 - Event sets do not keep timing information
- Result: rule-based failure prediction system containing a database of indicative eventsets

[Vilalta, Ma 2002]



Event-set Method



PFM - 19.6.2012 - 25



Hidden Semi-Markov Model Prediction

- Components of complex systems depend on each other
- Dependencies lead to error patterns



- Fault-tolerant systems:
 - Failures occur only under certain conditions
 - Failure-prone conditions can be identified by specific error patterns
- Use pattern recognition to identify symptomatic situations

[Salfner, Malek 2007; Salfner 2008]

PFM - 19.6.2012 - 26



Approach

- Standard tool for pattern recognition: Hidden Markov Models
- Identify symptomatic patterns
 - Algorithmically
 - From recorded training data
 - Machine learning
- Additional assumption:
 - Time between events is decisive (temporal sequence analysis)
 - Standard Hidden Markov Models need to be extended
 - Development of a Hidden Semi-Markov Model (HSMM)
- The approach incorporates both type and time-of-occurrence of error events







- Discrete Time Markov Chains (DTMC) consist of states (1...N) and transition probabilities p_{ij} between states
- In Hidden Markov Models (HMM) each state can generate a symbol A,B,C according to probability distribution b_i(o_k)
- Hidden semi-Markov models (HSMMs) replace transition probabilities p_{ij} by time-continuous cumulative probability distributions g_{ij}(t)

Machine Learning: Two Steps

1. Training: Fit model parameters to training data



2. Prediction: Processing of runtime measurements





Contents

- Introduction
- Variable Selection
- Online Failure Prediction Overview
- Four Online Failure Prediction Techniques
- Assessing Failure Predictors
- Taking Action
- Summary

Recall and other Metrics

contingency table	True failure	True success	Sum
Failure alarm	Correct alarm (TP)	False alarm (FP)	# Alarms
No warning	Missing alarm (FN)	Correct no-alarm (TN)	# No-Alarms
Sum	# Failures	# Successes	# Total

- Precision: fraction of correct alarms: $precision = \frac{correct alarms}{total \ \sharp \ of \ alarms}$
- Recall: fraction of predicted failures: r
- $\text{recall} = \frac{\text{correct alarms}}{\text{total } \sharp \text{ of failures}}$
- False positive rate (fpr): false positive rate = $\frac{\text{false positives}}{\ddagger \text{ of successes}}$
- True positive rate is equal to recall

Receiver Operating Characteristics (ROC)



- Plot true positive rate (recall) over false positive rate for various thresholds
- Threshold ∞ : tpr and fpr equal to zero
- Threshold $-\infty$: tpr and fpr equal to one



Precision-Recall-Plots



- Plot precision over recall for various thresholds
- Threshold ∞ : precision equal to one, recall equal to zero
- Threshold -∞ : precision equal to ratio of positive and negative examples, recall equal to one



Skalar Metrics

- ROC, Precision/Recall diagrams etc. are graphs
- Good for visual inspection, bad for algorithmic decisions
- Goal: obtain one real number to evaluate "quality" of predictor
- Examples:
 - Precision-Recall-Breakeven: Value at which precision and recall cross
 - Area-under-curve (AUC): Area under the ROC curve
 - F-Measure: Harmonic mean of precision and recall:

$$F-Measure = \frac{2 * precision * recall}{precision + recall}$$



Contents

- Introduction
- Variable Selection
- Online Failure Prediction Overview
- Four Online Failure Prediction Techniques
- Assessing Failure Predictors
- Taking Action
- Summary



Taking Action

- Failure prediction is only the first step in managing faults proactively
- After a potential failure has been predicted, an action must be taken in order to
 - Avoid the failure (prevent it from occurring)
 - Prepare the system such that TTR can be reduced
 - (See lecture seven)
- In general, the following steps have to be performed:



Taxonomy of Reaction Methods





Effects of Proactive Methods

- Downtime avoidance improves MTTF
- Downtime minimization reduces MTTR:



 However, In case of frequent false positive and false negative predictions, proactive fault management can also reduce availability!



Contents

- Introduction
- Variable Selection
- Online Failure Prediction Overview
- Four Online Failure Prediction Techniques
- Assessing Failure Predictors
- Taking Action
- Summary



Summary

- Dynamics and complexity of today's systems require adaptive and proactive mechanisms to handle faults
- Proactive Fault Management (PFM) builds on
 - Continuously observing the system
 - Predict whether a failure is coming up
 - In case of an upcoming failure:
 - o Analyze the fault that causes the upcoming failure
 - o Decide what to do: Either try to avoid the failure or prepare repair mechanisms for the upcoming failure
- Since online failure predictors build on monitoring data: The best set of variables need to be identified: Variable selection
- Analyses of PFM suggest that PFM has the potential to enhance system availability by up to an order of magnitude.



Thanks! Questions?