



Zentrum für sichere Informationstechnologie – Austria
Secure Information Technology Center – Austria

A-1040 Wien, Weyringergasse 35
 Tel.: (+43 1) 503 19 63–0
 Fax: (+43 1) 503 19 63–66

A-8010 Graz, Inffeldgasse 16a
 Tel.: (+43 316) 873-5514
 Fax: (+43 316) 873-5520

<http://www.a-sit.at>
 E-Mail: office@a-sit.at

EREIGNIS-KORRELATION
TECHNOLOGIEBEOBACHTUNG
VERSION 1.0 6/2/2009 10:00 PM

Peter Teufl – peter.teufl@a-sit.at

1 **Zusammenfassung:** Im A-SIT Projekt ***Ereignis-Korrelation*** wurde eine neuartige Methode
 2 entwickelt, um unterschiedliche Daten (Ereignisse) von unterschiedlichen Sensoren zu
 3 kombinieren und für weitere Analysen aufzubereiten. Dabei werden die Daten der einzelnen
 4 Sensoren und ihr gemeinsames Auftreten in einem assoziativen Netzwerk abgebildet. Durch die
 5 Anwendung von Spreading Activation werden Aktivierungsmuster erzeugt, die dann für weitere
 6 Analysen verwendet werden. Diese Arbeit beschreibt das Problem der Ereignis-Korrelation, stellt
 7 das neuartige Framework vor und evaluiert es anhand von Echtdateien.

Inhaltsverzeichnis

Inhaltsverzeichnis	1
Abbildungs- und Tabellenverzeichnis	2
1 Einleitung	3
2 Ereignis-Korrelation (Event Correlation)	5
3 Methoden	7
3.1 Semantische Netzwerke/Assoziative Netzwerke	7
3.2 Spreading Activation (SA)	7
4 Unüberwachtes Lernen	8
5 Überblick über das Framework	8
6 Framework im Detail	9
6.1 L1 – Feature Extraction	9
6.2 L2 – Node Generation	9
6.3 L3 – Network Generation	10
6.4 L4 – Generation of Activation Patterns	11
6.5 L5 – Analysis	11
7 Ergebnisse	13
7.1 Unüberwachtes Finden von Clustern	14
7.2 Relationen	15
7.3 Suche	15
7.4 Anomalieerkennung	16
8 Zusammenfassung und Ausblick	18
9 Referenzen	19
10 Anhänge	20
Historie	21

Abbildungs- und Tabellenverzeichnis

Abbildung 1 - Überblick über Ereignis-Korrelation	6
Abbildung 2 - Überblick über das vorgestellte Framework	12
Abbildung 3 - Aktivierungsenergie der Trainingsdaten (rot) und der Anomalien (schwarz)	17
Tabelle 1 - Beschreibung der Features des Automobiles Datensatzes	13
Tabelle 2 - Ergebnisse des unüberwachten Clusterings.....	15
Tabelle 3 - Relationen.....	15
Tabelle 4 - Suchergebnisse	16
Tabelle 5 - Aktivierungsenergie der Trainingsdaten und Anomaliedaten	17

1 Einleitung

Um der steigenden Anzahl von Angriffen auf IT Systeme entgegen zu wirken, werden Intrusion Detection Systeme (IDS, zu deutsch: Eindringlings-Erkennungs-Systeme) eingesetzt, die etwaige Angriffe oder Anomalien in einer IT Umgebung erkennen sollen. Dabei kann ein IDS Analysen auf unterschiedlichen Ebenen einer IT Umgebung durchführen. Diese Ebenen können von der Analyse des reinen Netzwerkverkehrs bis zur Analyse von Applikationen oder des Benutzerverhaltens auf einzelnen Maschinen reichen. Jede dieser Ebene spielt je nach Art des Angriffs eine unterschiedliche Rolle und kann besser oder schlechter für die Erkennung verwendet werden. Ebenso ist es möglich, dass sich ein Angriff nur durch die Analyse mehrerer unterschiedlichen Ebenen erkennen lässt. Die für das Sammeln der Daten zuständigen Komponenten werden als Sensoren bezeichnet.

Die Analyse der von den Sensoren gelieferten Daten ist dabei eine der wichtigsten Komponenten eines IDS. Folgende Punkte müssen dabei beachtet werden:

- **Analyse:** Es gibt Fälle wo einzelne Sensoren ausreichen, um einen Angriff zu erkennen. Dies ist aber nicht der Normalfall, und man muss davon ausgehen, dass erst die Kombination von unterschiedlichen Sensordaten eine Erkennung zulässt.
- **False Negatives:** Im Idealfall darf die Analyseeinheit des IDS keine Angriffe übersehen. Ist dies trotzdem der Fall, so spricht man von *False Negatives*, also Daten die fälschlicherweise als normal eingestuft wurden.
- **False Positives:** Im Gegensatz zu False Negatives beschreiben *False Positives* Ereignisse, die vom IDS fälschlicherweise als Angriff erkannt werden. Auch solche Ereignisse schränken die Qualität des IDS ein, da sie manuellen Analyseaufwand verursachen.

In dieser Arbeit konzentrieren wir uns auf die Analyse von Sensordaten (oder auch Ereignisse) die von Sensoren auf unterschiedlichen Ebenen erzeugt werden. Das Zusammenfassen und die Analyse dieser Ereignisse, die von unterschiedlichen Sensoren erzeugt werden, bezeichnet man als Ereignis-Korrelation.

Das vorgestellte Framework behandelt dabei folgende Punkte:

- **Unüberwachtes Lernen:** Wir gehen von einem System mit unterschiedlichen Sensoren aus. Jeder dieser Sensoren liefert Daten (Ereignisse), die Schlüsse über das zu analysierende System zulassen. Dabei gehen wir davon aus, dass wir kein Wissen über die Zusammenhänge dieser Ereignisse haben. Das Framework lernt dabei selbstständig die Relationen zwischen einzelnen Ereignissen und fasst das Gesamtverhalten nach Ähnlichkeit zusammen. Im Falle von Netzwerk IDS Daten kann dies dann die automatische Erkennung von unterschiedlichen Clustern sein: Denial of Service (DOS) Angriffe, Probes, normaler Verkehr.
- **Analyse von unterschiedlichen Sensordaten:** Unterschiedliche Sensordaten haben unterschiedliche Wertbereiche und können entweder symbolische oder kontinuierliche Daten liefern. Diese unterschiedlichen Daten erschweren die direkte unüberwachte Analyse. Aus diesem Grund fügt das Framework noch eine Zwischenschicht ein, die die Daten transformiert, sodass diese unüberwachte Analyse wieder möglich wird.
- **Korrelation der Sensordaten:** Die Betrachtung der unterschiedlichen Ereignisse und deren Relationen untereinander (z.B.: gemeinsames Auftreten) werden vom Framework erlernt und ermöglichen Schlüsse über die Zusammenhänge der einzelnen Ereignisse.

- 57
- 58
- 59
- 60
- 61
- **Anomalieerkennung:** Dabei wird für den normalen Verkehr ein Modell trainiert. Unbekannte Muster – in diesem Fall Angriffe - die während dem Training nicht vorhanden waren, sollen dann als Anomalie erkannt werden. Mit Hilfe der Aktivierungsenergie der im Framework vorgestellten Aktivierungsmuster kann diese Anomalieerkennung realisiert werden, da diese Energie bei Anomalien schwächer ist.

2 Ereignis-Korrelation (Event Correlation)

Das Sammeln von Sensordaten, die aufgrund von Ereignissen (Events) erzeugt wurden und deren Auswertung wird als Ereignis-Korrelation bezeichnet. Wenn wir dabei das Wort Korrelation verwenden, um die Verknüpfung unterschiedlicher Ereignisse zu bezeichnen, dann stellt sich hier gleich die Frage wie das Auftreten unterschiedlicher Ereignisse analysiert und klassifiziert werden kann. Hierbei können auf jeden Fall bestimmte Basisfunktionen angegeben werden, die von einem Ereignis-Korrelations-System (EKS) erfüllt werden müssen:

- **Filtern von Ereignisse:** Im Falle eines Angriffs werden die Sensoren eines IDS unterschiedliche Daten liefern, die von einem EKS korrekt interpretiert werden müssen. Ein EKS soll dabei die gemeinsame Ursache dieser Sensordaten erkennen und klassifizieren. Der Betreiber des EKS erhält nicht mehr die Datenströme der einzelnen Sensoren, sondern eine Zusammenfassung, die das aktuelle Verhalten auf höherer Ebene beschreibt.
- **Erkennen und Repräsentation von Verhaltensweisen:** Die Kernkomponente eines EKS basiert auf der Analyse der Sensordaten, der Interpretation dieser Analyse und der Repräsentation. Vor allem die Repräsentation spielt eine sehr wichtige Rolle, da es möglich sein muss, dem Benutzer auf möglichst verständliche Weise ein Bild des aktuellen Status zu liefern.

Es werden nun die Basiskomponenten eines EKS beschrieben:

Sensor: Ein Sensor liest beliebige Daten aus, wendet Vorverarbeitungsschritte auf diese Daten an und übergibt Ereignisse an höhere Schichten des EKS, wenn bestimmte Bedingungen erfüllt sind. Ein einfaches Beispiel für so eine Bedingung ist ein Temperatur-Sensor der Ereignisse nur an höhere Schichten weiterleitet, wenn die gemessene Temperatur größer als ein bestimmter Schwellwert ist.

Ereignis Modell: Dieses Modell beschreibt die Relationen zwischen unterschiedlichen Ereignissen. Dabei kann es sich um ein einfaches, auf Regeln basierendes Modell oder um ein komplexes mathematisches Modell handeln. Für die Erzeugung der Modelle können unterschiedliche Verfahren verwendet werden – z.B. Algorithmen aus dem maschinellen Lernen.

Korrelations-Komponente: Diese Komponente vergleicht zu überprüfende Ereignisse mit den trainierten Modellen und zieht daraus Schlüsse über das zu analysierende Verhaltensmuster.

Bei der Korrelationskomponente spielen dabei folgende Punkte eine Rolle:

Selektion der Informationen: Es kann sein, dass das gleiche Ereignis von verschiedenen Sensoren gleichzeitig gemeldet wird. In diesem Fall kann es sinnvoll sein, mehrfach vorhandene Ereignisse zu löschen. Allerdings muss darauf geachtet werden, dass in manchen Fällen auch die Information des wiederholten Auftretens eines Ereignisses für die Erkennung von Verhaltensmustern verwendet werden kann.

Filtern: Das EKS muss in der Lage sein, aus den gesamten Ereignissen jene zu extrahieren, die für die weitere Analyse des Verhaltens den größten Informationsgehalt bringen.

110 **Ereignis Repräsentation:** Abhängig vom Zustand des EKS kann es Sinn machen, Ereignisse
111 zu unterdrücken oder in einer anderen Art zu repräsentieren.

112
113 **Zählen von Ereignissen:** EKS können diese sehr einfache Methode anwenden, um einen
114 Alarm zu generieren, wenn die Anzahl bestimmter Ereignisse einen bestimmten Schwellwert
115 überschreitet. In dieser Arbeit konzentrieren wir uns aber auf viel höher entwickelte Analyse-
116 Werkzeuge.

117
118 **Gewichtung von Ereignissen:** Das EKS hat die Aufgabe, unterschiedliche Ereignisse je nach
119 Informationsgehalt zu gewichten.

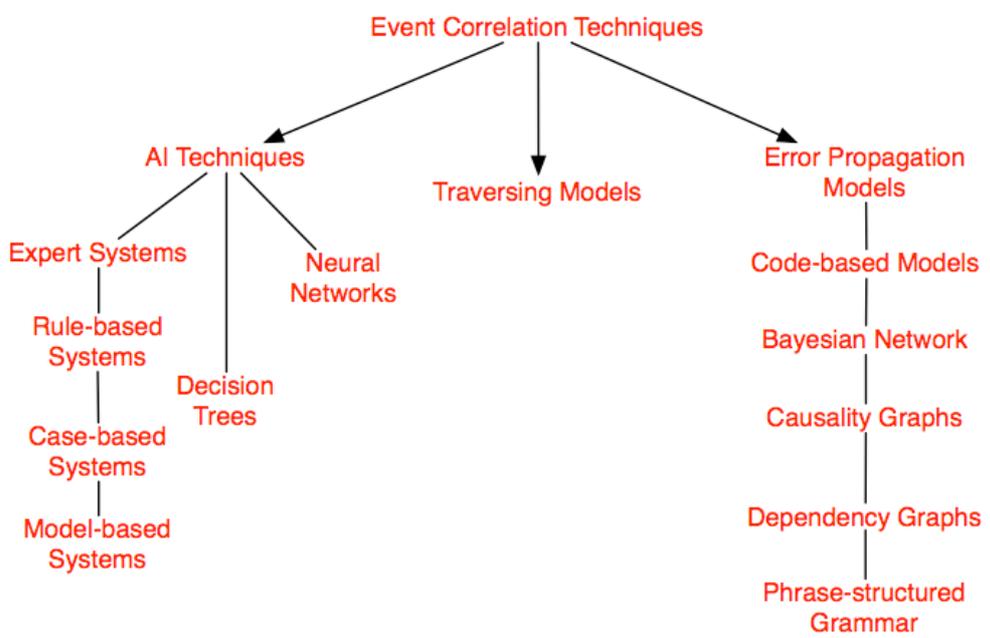
120
121 **Generalisierung von Ereignissen:** Das EKS soll in der Lage sein, spezifische Ereignisse einer
122 allgemeineren Klasse zuzuordnen. Beispiele dafür sind z.B. das Erkennen eines speziellen DOS
123 Angriffs, der zur allgemeinen Klasse der DOS Angriffe zugeordnet wird.

124 **Zeitliche Zusammenhänge:** Die zeitliche Abfolge von Ereignissen kann Schlüsse auf das
125 darunter liegende Verhalten zulassen und soll von einem EKS berücksichtigt werden.

126
127 **Clustern von Ereignissen:** Das Clustern oder Zusammenfassen von ähnlichen
128 Verhaltensmustern spielt vor allem dann eine Rolle wenn noch keine Informationen über die
129 möglichen Muster vorliegen. Das System versucht dann selbständig ähnliche Muster zu finden
130 und in Clustern zusammenzufassen (z.B. unterschiedliche Angriffe, normaler Verkehr).

131 **Modelle basierend auf Künstlicher Intelligenz (KI):** Solche Modelle versuchen, menschliche
132 Problemlösungs-Strategien abzubilden und für die Analyse von Ereignissen zu verwenden. Es
133 können dabei Expertensysteme oder eine Vielzahl anderer Methoden aus der Künstlichen
134 Intelligenz angewendet werden.

135



136
137

Abbildung 1 - Überblick über Ereignis-Korrelation

3 Methoden

138 Dieser Abschnitt gibt einen kurzen Überblick über die Algorithmen und Methoden, die vom
139 Framework verwendet werden.

3.1 Semantische Netzwerke/Assoziative Netzwerke

140 Semantische Netzwerke [1] bestehen aus Knoten, die Informationen repräsentieren. Diese
141 Knoten werden untereinander mit Kanten verbunden, die die Relationen zwischen den
142 Knoten beschreiben. Dabei werden je nach Relation unterschiedliche Kantentypen eingesetzt.
143 Gute Beispiele für semantische Netzwerke sind WordNet [2], [3] und GermaNet [1].
144 Assoziative Netzwerke sind allgemeiner als semantische Netzwerke und verwenden
145 generische Kanten ohne Labels. Diese Kanten können auch gewichtet werden und sind dann
146 in der Lage, etwas über die Stärke der Relation auszusagen. In dieser Arbeit werden
147 assoziative Netzwerke eingesetzt, um die Relationen zwischen Ereignissen abzubilden.

3.2 Spreading Activation (SA)

148 Spreading Activation Algorithmen (SA Algorithmen) werden verwendet, um Informationen
149 aus semantischen/assoziativen Netzwerken zu extrahieren. Dabei werden zuerst ein oder
150 mehrere Knoten im Netzwerk aktiviert. In den folgenden vom Algorithmus durchgeführten
151 Iterationen wird diese Aktivierung je nach Stärke der Relationen auf benachbarte Knoten
152 weitergegeben. Die Anzahl der Iterationen bestimmt die Tiefe der berücksichtigten
153 Nachbarschaftsbeziehungen.

154 In dieser Arbeit werden SA Algorithmen für folgende Komponenten eingesetzt:

155

- 156 • **Relationen zwischen Ereignissen:** Durch die Aktivierung von Knoten, die Ereignisse
157 repräsentieren, und die anschließende Anwendung von SA Algorithmen kann
158 festgestellt werden, welche anderen Ereignisse mit den aktivierten Ereignissen
159 assoziiert sind und wie stark diese Assoziationen sind.
- 160 • **Generierung der Aktivierungsmuster:** Für die weitere Analyse verwenden wir
161 sogenannte Aktivierungsmuster, die Auskunft über die im assoziativen Netzwerk
162 aktivierten Bereiche geben. Ein Aktivierungsmuster enthält alle Aktivierungswerte, des
163 assoziativen Netzwerks. Die Anzahl der Einträge entspricht somit der Anzahl der
164 Knoten im Netzwerk.

4 Unüberwachtes Lernen

165 Unüberwachte Lernalgorithmen analysieren Daten und erkennen Ähnlichkeiten innerhalb
166 dieser Daten. Ähnliche Daten können dann in Gruppen/Cluster zusammengefasst werden.
167 Speziell im Bereich Ereignis-Korrelation bei IDS spielt diese Lernmethode eine sehr wichtige
168 Rolle, da man normalerweise nicht über Daten von Angriffen verfügt, die für das Training
169 verwendet werden können. Beispiele für solche Lernalgorithmen sind Self Organizing Maps
170 [6], Neural Gas und dessen Derivate [4], [5] und der EM Algorithmus [7].

5 Überblick über das Framework

171 Das vorgestellte Framework besteht aus fünf Schichten, die in Abbildung 2 (Seite 12)
172 dargestellt werden.

- 173 • **L1 – Feature Extraction:** In diesem Schritt werden manuell die kontinuierlichen und
174 die symbolischen Features extrahiert. Kontinuierliche Werte können zu größeren
175 Gruppen zusammengefasst werden, wenn es die Wertbereiche erlauben. Das Erstellen
176 von Gruppen verbessert die Performance ist aber keine Notwendigkeit.
- 177 • **L2 – Erstellen der Knoten:** In diesem Schritt werden die Knoten des assoziativen
178 Netzwerks erstellt. Dabei wird zwischen den unterschiedlichen Features
179 unterschieden:
 - 180 ○ **Kontinuierliche Werte:** Bevor diese Werte in das Netzwerk integriert werden
181 können, muss ein unüberwachter Lernalgorithmus angewendet werden. Die
182 gefunden Cluster werden dann für die Knoten im Netzwerk verwendet.
 - 183 ○ **Symbolische Werte:** Diese Werte können direkt als Knoten ins Netzwerk
184 integriert werden.
- 185 • **L3 – Erstellen des Netzwerks:** Die Relationen im assoziativen Netzwerk werden
186 aufgrund des gemeinsamen Auftretens unterschiedlichen Ereignisse erstellt. Die Stärke
187 oder das Gewicht der Relationen hängt dabei von der Anzahl des gemeinsamen
188 Auftretens dieser Ereignisse ab.
- 189 • **L4 – Bestimmen der Aktivierungsmuster:** Die Aktivierungsmuster werden nun mit
190 Hilfe des assoziativen Netzwerks und einem SA Algorithmus erstellt. Dazu werden die
191 Knoten der auftretenden Ereignisse im assoziativen Netzwerk aktiviert und
192 anschließend der SA Algorithmus angewendet, um die Aktivierung auf benachbarte
193 Knoten zu verteilen. Anschließend werden die Aktivierungswerte aller Knoten im
194 Netzwerk extrahiert und in einem Vektor gespeichert. Dieser Vektor entspricht dem
195 Aktivierungsmuster und kann für die weitere Analyse verwendet werden.
- 196 • **L5 – Analyse:** Die erstellten Aktivierungsmuster können nun für die weitere Analyse
197 verwendet werden. Dabei können sowohl überwachte als auch unüberwachte Analyse-
198 verfahren angewendet werden. Außerdem ist es mit Hilfe von SA Algorithmen möglich,
199 die Relationen zwischen Ereignissen zu analysieren.

6 Framework im Detail

200 In diesem Abschnitt werden die einzelnen Schichten des Frameworks detailliert beschrieben.

6.1 L1 – Feature Extraction

201 Dieser Layer behandelt die Extraktion der Daten, die von unterschiedlichen Sensoren zur
202 Verfügung gestellt werden. Diese Daten können in zwei Kategorien eingeteilt werden:

203

- 204 • **Daten für die ein Distanzmaß definiert werden kann:** In diesem Fall ist es möglich,
205 die Distanz zwischen zwei unterschiedlichen Messwerten zu definieren. Ein Beispiel
206 dafür wäre ein Feature, das die Anzahl der verlorenen Pakete pro Verbindung angibt.
207 Durch die Existenz des Distanzmaßes ist es nun möglich, die Ähnlichkeiten zwischen
208 unterschiedlichen Fehlerraten zu berechnen – eine Fehlerrate von 10% ist näher zu
209 einer Fehlerrate von 12% als zu einer Fehlerrate von 80%.
- 210 • **Daten für die kein Distanzmaß definiert werden kann:** In diesem Fall handelt es
211 sich um symbolische Werte, die nicht über eine Distanz in Relation gesetzt werden
212 können. Beispiele für solche Werte sind Protokoll Typen wie TCP, UDP, ICMP. Es ist
213 nicht möglich, Aussagen der Art “TCP ist näher zu UDP als zu ICMP” zu treffen, da kein
214 Distanzmaß zwischen den Protokolltypen definiert ist.

215 Diese Einteilung muss für alle Features getroffen werden, da beide Kategorien unterschiedlich
216 behandelt werden müssen. Nach dem Zuordnen der Features können für die Werte mit
217 Distanzmaß Gruppen gebildet werden, die Features mit ähnlichen Wertbereichen
218 kombinieren. Beispiele dafür sind Prozentwerte, die Auskunft über unterschiedliche
219 Fehlerraten geben. Die Gründe für diese Gruppenbildung werden in L2 besprochen.

6.2 L2 – Node Generation

220 Dieser Layer ist für die Erstellung der Knoten im assoziativen Netzwerk verantwortlich. Dazu
221 muss nun eine Unterscheidung zwischen den in L1 festgelegten Kategorien getroffen werden:
222 Symbolische Daten (z.B. UDP, TCP, ICMP) können direkt als Knoten in das assoziative
223 Netzwerk integriert werden. Bei den Daten für die es möglich ist ein Distanzmaß zu
224 definieren, muss allerdings noch ein weiterer Verarbeitungsschritt eingefügt werden. Eine
225 direkte Repräsentation dieser Werte als Knoten ist nicht möglich, da sie typischerweise nicht
226 abzählbar sind. Selbst wenn dies der Fall ist, macht es keinen Sinn, jeden möglichen Wert als
227 Knoten zu repräsentieren, da dann die Distanz-Information verloren geht. Diese Information
228 sagt etwas über die Ähnlichkeit zwischen einzelnen Werten aus und kann daher verwendet
229 werden um Häufungen (Cluster) von ähnlichen Werten zu finden. Diese Cluster
230 repräsentieren dann die Knoten im assoziativen Netzwerk.

231 Um diese Cluster zu finden, sind folgende Schritte notwendig:

232

- 233 1. Jede Feature Gruppe, bestehend aus einem oder mehreren Features, wird mit einem
234 unüberwachten Lernalgorithmus analysiert. Dieser Lernalgorithmus findet Cluster
235 innerhalb der Daten. In dieser Arbeit wird hierfür ein Algorithmus verwendet, der auf
236 Growing Neural Gas basiert.
- 237 2. Die gefundenen Cluster stellen die Knoten dar, die im assoziativen Netzwerk
238 hinzugefügt werden.
- 239 3. Werte können nun wie folgt auf die im Netzwerk vorhandenen Knoten abgebildet
240 werden: Für einen unbekanntem Wert wird der am nächsten liegende Cluster
241 bestimmt. Da die Cluster direkt als Knoten repräsentiert werden, stellt der gefundene
242 Cluster auch direkt den im Netzwerk gefundenen Knoten dar.

243 Die in L1 besprochene Gruppenbildung macht aus Sicht der Komplexität Sinn, da der
244 unüberwachte Lernalgorithmus nur einmal pro Gruppe angewendet werden muss.

6.3 L3 – Network Generation

245 Es werden die zu analysierenden Daten nun dazu verwendet, die Relationen zwischen den
246 Ereignissen im assoziativen Netzwerk herzustellen. Dazu müssen die Daten in einem Format
247 vorliegen, das das gemeinsame Auftreten von Ereignissen (Sensordaten) beschreibt. Die
248 genaue Definition des gemeinsamen Auftretens hängt nun von der Art der zu analysierenden
249 Daten ab: z.B. ein gewisser Zeitraum oder eine andere Abhängigkeit, die nicht auf Zeit basiert
250 (z.B.: im Bereich Intrusion Detection: Ereignisse mit gleicher Quelladresse).

251 Sensordaten die im Rahmen dieser definierten Gemeinsamkeit aufgezeichnet wurden, stellen
252 einen Datenvektor dar, der für das Training des assoziativen Netzwerks verwendet wird.
253 Dabei wird für jeden Datenvektor in den Trainingsdaten wie folgt vorgegangen:

254

- 255 1. Die Features (Sensordaten, Ereignisse) pro Datenvektor werden anhand der in L1
256 getroffenen Einteilung und Gruppierung analysiert. Features, für die kein Distanzmaß
257 festgelegt werden kann, werden direkt auf die Knoten im Netzwerk abgebildet. Für die
258 anderen Features werden die in L2 trainierten unüberwachten Modelle hergenommen
259 und die am nächsten liegenden Cluster bestimmt. Diese Cluster repräsentieren die
260 Knoten im Netzwerk.
- 261 2. Es werden nun Links zwischen allen Ereignissen im Datenvektor hergestellt. Ist ein
262 Link zwischen zwei Ereignissen noch nicht vorhanden, so wird dieser mit dem
263 Gewicht 1 initialisiert. Wenn der Link schon vorhanden ist, dann wird dieses Gewicht
264 um 1 erhöht.

265 Nach der Analyse aller Datenvektoren hat man nun ein assoziatives Netzwerk, das die
266 Assoziationen zwischen Ereignissen über die Kanten zwischen den Knoten herstellt. Das
267 Gewicht der Kanten entspricht der Anzahl des gemeinsamen Auftretens. Die Knoten selbst
268 repräsentieren die Werte der unterschiedlichen Features.

269 Um in den weiteren Layers den SA Algorithmus anwenden zu können, müssen diese Gewichte
270 normiert werden, so dass Maximalwerte von 1.0 vorkommen. Es gibt unterschiedliche
271 Strategien, um diese Normierung durchzuführen:

272

- 273 • **Globale Maximum Norm:** In diesem Fall wird das maximale Gewicht im Netzwerk
274 genommen und alle anderen Gewichte werden damit normiert.
- 275 • **Lokale Maximum Norm:** In diesem Fall wird das maximale Gewicht aller Kanten, die
276 von einem Knoten ausgehen, genommen und alle Kanten dieses Knotens mit diesem
277 Gewicht normiert. Dies wird für alle Knoten durchgeführt.
- 278 • **Lokale Summen Norm:** In diesem Fall wird die Summe der von einem Knoten
279 ausgehenden Gewichte genommen und jedes Gewicht mit dieser Summe normiert.

280 Bei der Anwendung einer lokalen Norm geht im Gegensatz zu einer globalen Normierung die
281 Symmetrie der Kanten verloren.

282 Für die in dieser Arbeit vorgestellten Ergebnisse wurde die lokale Summennorm verwendet.
283 Die Gründe dafür sind:

- 284 • Eine lokale Normierung ermöglicht es, Assoziationen zu verstärken, die zwar global
285 gesehen schwach sind, aber lokal eine wichtige Bedeutung haben.
- 286 • Die Summennorm bewirkt, dass für jeden Knoten eine Aktivierungsenergie von 1.0 zur
287 Verfügung steht. D.h. die Aktivierungsenergie wird auf alle ausgehenden Kanten
288 verteilt. Je mehr Kanten vorhanden sind, umso weniger Energie steht für jede einzelne

289 Kante zur Verfügung. Dies entspricht bereits der Funktion eines einfachen Fanout
290 Faktors.

291 Nach der Normierung steht nun ein assoziatives Netzwerk zur Verfügung, das für die weiteren
292 Analysen verwendet wird.

6.4 L4 – Generation of Activation Patterns

293 Die für das Training des assoziativen Netzwerks verwendeten Datenvektoren werden nun
294 verwendet, um die sogenannten Aktivierungsmuster zu generieren. Dazu wird für jeden
295 Datenvektor wie folgt vorgegangen:

- 296 1. Es werden die Features aus dem Datenvektor extrahiert und auf die Knoten im
297 Netzwerk abgebildet. Dabei wird gleich wie in L3 vorgegangen.
- 298 2. Alle Knoten erhalten den Aktivierungswert 1.0.
- 299 3. Durch Anwenden des SA Algorithmus werden die Aktivierungswerte auf die Nachbarn
300 verteilt. Die Stärke der Aktivierung der Nachbarn hängt dabei von dem Gewicht der
301 Kanten ab. In dieser Arbeit wird nur eine Iteration verwendet, da sonst die Aktivierung
302 zu weit im Netzwerk verbreitet wird.

303 Nach Anwendung des SA Algorithmus auf die im Netzwerk aktivierten Knoten werden die
304 Aktivierungswerte aller Knoten im Netzwerk ausgelesen. Diese Werte werden dann in einen
305 Vektor eingetragen. Dieser Vektor wird als Aktivierungsmuster bezeichnet.

6.5 L5 – Analysis

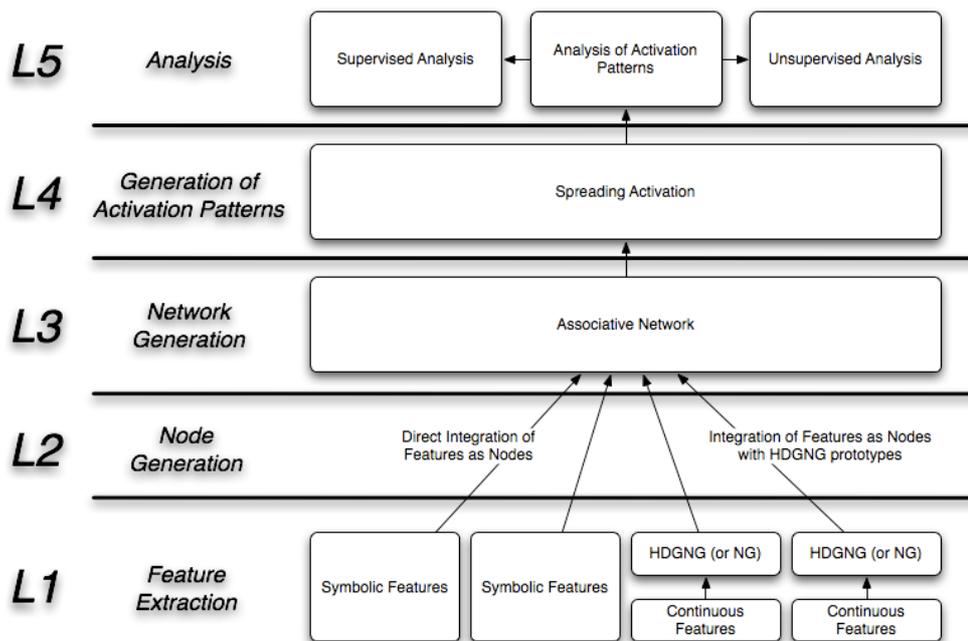
306 In L4 wurden die Aktivierungsmuster für alle Datenvektoren erstellt. Diese
307 Aktivierungsmuster können nun für weitere Analyse Schritte verwendet werden. Dazu
308 gehören:

- 309 • **Unüberwachtes Finden von Clustern:** Es kann ein beliebiger unüberwachter
310 Lernalgorithmus eingesetzt werden, der Cluster innerhalb der Aktivierungsmuster
311 findet. Solche Cluster repräsentieren unterschiedliche Vorkommnisse. Im Intrusion
312 Detection Bereich können dies z.B. unterschiedliche Typen von Angriffen sein (DDOS,
313 Syn Flood).
- 314 • **Überwachtes Lernen:** Wenn man bereits Informationen über die Zugehörigkeit (z.B.
315 bekannte Angriffe) der zu analysierenden Datenvektoren hat, können überwachte
316 Lernalgorithmen eingesetzt werden, um ein Modell für die Daten zu lernen. Dieses
317 Modell kann dann benutzt werden, um unbekannte Datenvektoren zu klassifizieren.
- 318 • **Suche nach assoziierten Ereignissen:** Das assoziative Netzwerk und der SA
319 Algorithmus können dazu verwendet werden, um verwandte Ereignisse zu finden.
320 Dazu werden die Knoten eines oder mehrerer Ereignisse im Netzwerk aktiviert.
321 Danach wird der SA Algorithmus verwendet, um die Aktivierung auf Nachbarknoten zu
322 verteilen. Die Stärke der Aktivierung dieser Nachbarknoten lässt den Schluss über die
323 Stärke der Assoziation zu anderen Ereignissen zu.
- 324 • **Vergleich von Aktivierungsmustern:** Die Stärke der Aktivierung der einzelnen
325 Knoten wird im Aktivierungsmuster gespeichert. Der Vergleich zwischen zwei
326 Aktivierungsmustern gibt nun Auskunft darüber, welche Bereiche im Netzwerk (Typen
327 von Ereignissen) aktiv sind und wo die Unterschiede liegen.
- 328 • **Anomalieerkennung:** Als Aktivierungsenergie bezeichnen wir die Summe der
329 Aktivierungswerte der Knoten des assoziativen Netzwerks nach Anwendung des SA
330 Algorithmus. Die Aktivierungsenergie gibt Auskunft darüber, wie gut vorgegebene
331 Muster in das trainierte Netzwerk passen. Muster, die ähnlich zu den im Training
332 vorhanden Muster sind, ergeben eine höhere Aktivierungsenergie, da die Links
333 zwischen den einzelnen Features stärker sind. Bei unbekanntem Mustern ist es

334
335
336
337

typischerweise der Fall, dass Werte der einzelnen Features in anderen Kombinationen auftreten und somit Knoten im Netzwerk aktivieren, deren Verknüpfung untereinander schwächer ist. Dies spiegelt sich in der Aktivierungsenergie wider und kann für die Anomalieerkennung genutzt werden.

Feature Extraction



338

Abbildung 2 - Überblick über das vorgestellte Framework

7 Ergebnisse

339 Um die Möglichkeiten des Frameworks zu zeigen, wird der *Automobiles* Datensatz aus dem
340 UCI Machine Learning Repository [9] genommen. In diesem Datensatz werden Autos von
341 unterschiedlichen Herstellern aus dem Jahr 1985 beschrieben. Für diese Beschreibung
342 werden 25 symbolische und Distanz-basierte Features genommen. Tabelle 1 beschreibt die
343 Bedeutung der Features und deren möglichen Werte.
344

ID	Feature	Werte
1	symboling	-3, -2, -1, 0, 1, 2, 3
2	normalized-losses	continuous from 65 to 256.
3	make	alfa-romero, audi, bmw, chevrolet, dodge, honda, isuzu, jaguar, mazda, mercedes-benz, mercury, mitsubishi, nissan, peugot, plymouth, porsche, renault, saab, subaru, toyota, volkswagen, volvo
4	fuel-type	diesel, gas
5	aspiration	std, turbo
6	num-of-doors	four, two
7	body-style	hardtop, wagon, sedan, hatchback, convertible
8	drive-wheels	4wd, fwd, rwd
9	engine-location	front, rear
10	wheel-base	continuous from 86.6 to 120.9.
11	length	continuous from 141.1 to 208.1.
12	width	continuous from 60.3 to 72.3.
13	height	continuous from 47.8 to 59.8.
14	curb-weight	continuous from 1488 to 4066.
15	engine-type	dohc, dohcv, l, ohc, ohcf, ohcv, rotor.
16	num-of-cylinders	eight, five, four, six, three, twelve, two.
17	engine-size	continuous from 61 to 326.
18	fuel-system	1bbl, 2bbl, 4bbl, idi, mfi, mpfi, spdi, spfi.
19	bore	continuous from 2.54 to 3.94.
20	stroke	continuous from 2.07 to 4.17.
21	compression-ratio	continuous from 7 to 23.
22	horsepower	continuous from 48 to 288.
23	peak-rpm	continuous from 4150 to 6600.
24	city-mpg	continuous from 13 to 49.
25	highway-mpg	continuous from 16 to 54.

Tabelle 1 - Beschreibung der Features des Automobiles Datensatzes

345
346 Warum wurde dieser Datensatz gewählt und wie ist die Verknüpfung zum Thema Ereignis
347 Korrelation?
348

- 349 • **Sensoren:** Die einzelnen Features können als unterschiedliche Sensoren betrachtet
350 werden, wobei von jedem Sensor unterschiedliche Aspekte abgedeckt werden. Dabei
351 gibt es sowohl Sensoren, die symbolischen Werte (vergleichbar mit Log-Einträgen von
352 IDS Systemen) und distanz-basierte Werte abdecken.
- 353 • **Ereignisse:** Die Werte, die von den einzelnen Sensoren geliefert werden, stellen die
354 Ereignisse dar, die wir aus dem Bereich Intrusion Detection kennen. Die Verknüpfung

355 dieser Ereignisse liefert Informationen über die Art des Fahrzeugs (bei IDS: z.B.:
356 Angriff) und ermöglicht Schlussfolgerungen über die Relationen der Ereignisse
357 untereinander (z.B. Motorstärke und Verbrauch, bei IDS: Protokoll-Typ und Anzahl der
358 Verbindungen pro Minute).

- 359 • **Verständlichkeit:** Die Analyse-Ergebnisse des Algorithmus können leicht überprüft
360 werden, da man mit vielen der Features des Datensatzes gut vertraut ist: z.B.:
361 Länge/Breite des Fahrzeugs, Stärke des Motors, Verbrauch Stadt/Überland usw.
- 362 • **Anzahl der Daten:** Im Datensatz werden 205 Fahrzeuge beschrieben. Diese
363 überschaubare Menge ermöglicht es, die Möglichkeiten des Frameworks besser zu
364 zeigen.

365 Für die Evaluierung werden nun vier Analysen durchgeführt:

- 366 • **Unüberwachtes Clustern:** Dabei werden ähnliche Fahrzeuge in Clustern/Gruppen
367 unüberwacht zusammengefasst.
- 368 • **Relationen zwischen den Features:** Dabei werden die Relationen der einzelnen
369 Features betrachtet und Fragen der Art "Wie hängen die Stärke des Motors und der
370 Verbrauch zusammen?" beantwortet.
- 371 • **Suche nach verwandten Mustern:** Dabei wird ein Suchmuster angegeben und nach
372 verwandten Mustern gesucht. Dabei kann es sich bei dem Suchmuster um ein Fahrzeug
373 aus dem Datensatz handeln, bei dem die Werte aller 25 Features festgelegt sind.
374 Weiters ist es auch möglich ausgewählte Features zu bestimmen und nach verwandten
375 Mustern zu suchen: z.B.: Festlegung von Verbrauch und Länge.

376 Für die Analyse der Daten wurden zwar alle Features verwendet, um aber die Relevanz der
377 Cluster zu zeigen, beschränkt sich die Tabelle dabei aber auf ein paar wenige leicht
378 verständliche Features: Marke, Länge, PS, MPG1 (Miles Per Gallon Stadt), MPG2 (Miles Per
379 Gallon Überland), Zylinder, Antrieb (Vorderrad, Hinterrad, Vierrad), Treibstoff
380 (Diesel/Benzin) und max. RPM.

7.1 Unüberwachtes Finden von Clustern

381 Bei dieser Analyse werden die Operationen der Layer 1-4 auf den Datensatz angewendet. Die
382 gewonnen Aktivierungsmuster werden anschließend mit einem unüberwachten
383 Lernalgorithmus analysiert. Dabei werden ähnliche Aktivierungsmuster in Gruppen (Cluster)
384 zusammengefasst. Die unüberwachte Analyse spielt hier eine wichtige Rolle, da wir davon
385 ausgehen, dass wir keine weiteren Informationen über die zu analysierenden Daten haben.
386 Bei dem vorliegenden Datensatz gibt es aber den Vorteil, dass wir über die einzelnen Features
387 und deren Relationen Bescheid wissen und somit eine bestimmte Clusterbildung erwarten
388 können. In diesem Fall darf man erwarten, dass Autos anhand ihrer Stärke und Größe
389 gruppiert werden. Dies lässt sich damit erklären, dass viele der vorhandenen Features diese
390 Eigenschaften beschreiben und untereinander korreliert sind. Ein Beispiel dafür wäre ein
391 Auto mit einem hohen Treibstoffverbrauch. In diesem Fall wird mit ziemlicher Sicherheit auch
392 ein starker Motor verwendet werden. Dies reflektiert sich z.B. auch in den Features Leistung
393 und Zylinder. Die Ergebnisse der unüberwachten Clusterbildung sind in Tabelle 2 zu sehen.

394

ID (Anzahl)	Marke	Länge	PS	MPG1	MPG2	Zyl.	Antrieb	B/D	RPM
C1 (71)	Toyota	188	92	19	24	4	RWD	Gas	5364
C2 (52)	Toyota	167	69	39	37	4	FWD	Gas	5364
C3 (9)	Toyota	176	69	30	32	4	FWD	Gas	4874
C4 (4)	Mazda	167	92	16	24	4	RWD	Gas	5916
C5 (24)	Toyota	172	92	23	32	4	FWD	Gas	5364
C6 (8)	Toyota	167	69	30	32	4	FWD	Gas	5364

C7 (13)	Toyota	176	92	23	32	4	FWD	Gas	4874
C8 (8)	Mercedes	188	180	16	24	4	RWD	Gas	4874
C9 (10)	Subaru	172	92	23	32	4	FWD	Gas	4874
C10 (6)	Toyota	176	114	23	28	4	FWD	Gas	4874

Tabelle 2 - Ergebnisse des unüberwachten Clusterings

395

396

397

398

399

400

401

Für die Tabelle wird für jeden Cluster der Mittelwert aller Aktivierungsmuster pro Cluster berechnet und ausgegeben. Da Toyota einen großen Bereich von unterschiedlichen Fahrzeugen abdeckt, ist diese Marke bei den Mittelwerten oft aktiviert. In der Tabelle kann man aber anhand des Verbrauchs und der PS sehen, dass die Fahrzeuge hauptsächlich nach diesen Kriterien gruppiert werden. Um einen guten Überblick über die 10 Cluster zu bekommen, wird auf die Liste der Gesamtergebnisse (*cluster.csv*) verwiesen.

7.2 Relationen

402

403

404

405

406

Bei dieser Analyse stellt man sich die Frage wie die Relationen zwischen unterschiedlichen Features sind. Ein Beispiel dafür ist die Frage "Mit welchen Werten für Motor Leistung, Verbrauch ist die Marke Porsche assoziiert". Es werden nun einige Beispiele in Tabelle 3 gezeigt:

	Marke	Länge	PS	MPG1	MPG2	Zyl.	Antrieb	B/D	RPM
Marke: Porsche <i>test-b.csv</i>	Porsche	167	205	16	24	6	RWD	Benz.	5961
Marke: VW <i>test-c.csv</i>	VW	172	92	27	32	4	FWD	Benz.	5364
MPG2: 40 <i>test-d.csv</i>	Chevrolet	157	70	37	42	4	FWD	Benz.	5364
PS: 200 <i>test-g.csv</i>	Porsche	167	205	16	24	6	RWD	Benz.	5961

Tabelle 3 - Relationen

407

408

409

410

411

412

413

414

415

Die in der Tabelle dargestellten Ergebnisse entsprechen den Erwartungen. Die Marke Porsche ist mit einer großen PS Zahl, einem Hinterradantrieb und einem hohen Verbrauch assoziiert. VW entspricht hier erwartungsgemäß genau dem Gegenteil. Das dritte Beispiel aktiviert einen Knoten der einem niedrigen Verbrauch entspricht. Hier sieht man, dass die Marke Chevrolet und ein schwacher Motor damit assoziiert sind. Auch das vierte Beispiel bringt die erwarteten Ergebnisse. Eine hohe PS Anzahl ist mit Porsche und einem hohen Verbrauch assoziiert. Diese Ergebnisse können leicht durch eine manuelle Überprüfung des Datensatzes verifiziert werden.

7.3 Suche

416

417

418

419

420

421

422

423

Hier soll ein Aktivierungsmuster vorgegeben werden und nach verwandten Mustern gesucht werden. Es kann sich hierbei um bestehende Aktivierungsmuster handeln, die im Zuge des Analyseprozesses der L1-L4 erzeugt wurden (z.B. Verwendung eines Autos aus dem Datensatz) oder um ein neues Aktivierungsmuster handeln, das für die Suche erstellt wird. Im zweiten Fall ist es auch möglich, nur wenige Features festzulegen, daraus ein Aktivierungsmuster zu erstellen und anschließend nach verwandten Mustern zu suchen. In Tabelle 4 werden nun Beispiele für beide Arten der Suche gezeigt:

	Marke	Länge	PS	MPG1	MPG2	Zyl.	Antrieb	B/D	RPM
M1	VW	172	58	37	48	4	FWD	Diesel	4874

test-e.csv									
4	Toyota	172	58	37	49	4	FWD	Diesel	4319
6	VW	167	92	27	32	4	FWD	Benz.	5364
11	Mazda	177	70	37	42	4	FWD	Diesel	4874
204	Jaguar	188	260	16	18	12	RWD	Benz.	4874
205	Mercedes	200	180	16	18	8	RWD	Benz.	4319
M2 test-f.csv	Jaguar	200	180	16	18	6	RWD	Benz.	4874
3	BMW	200	180	16	18	6	RWD	Benz.	5364
11	Peugeot	200	92	19	24	4	RWD	Benz.	4874
15	Audi	188	152	16	18	5	FWD	Benz.	5364
204	Honda	147	70	30	37	4	FWD	Benz.	5961
205	Chevrolet	147	58	46	48	3	FWD	Benz.	4874
M3 test-b.csv	Porsche	-	-	-	-	-	-	-	-
1	Porsche	167	205	16	24	6	RWD	Benz.	5961
6	Nissan	172	205	16	24	6	RWD	Benz.	5364
8	Mazda	172	205	16	24	6	RWD	Benz.	5364
203	VW	172	70	37	42	4	FWD	Diesel	4319
205	Toyota	176	70	30	32	4	FWD	Diesel	4874
M4 test-c.csv	VW	-	-	-	-	-	-	-	-
1	VW	172	70	37	42	4	FWD	Diesel	4319
13	Mazda	176	70	37	42	4	FWD	Diesel	4874
14	Toyota	167	58	37	42	4	FWD	Diesel	4874
203	Porsche	167	205	16	24	6	RWD	Benz.	5961
205	Jaguar	200	180	16	18	6	RWD	Benz.	4874
M5 test-a.csv	-	170	100	-	40	-	-	-	-
1	VW	172	70	37	42	4	FWD	Diesel	4319
2	Subaru	172	92	26	32	4	FWD	Benz.	4874
12	Nissan	172	92	26	32	4	FWD	Benz.	5364
203	Jaguar	188	260	16	18	12	RWD	Benz.	4874
205	Mercedes	200	180	16	18	8	RWD	Benz.	4319

Tabelle 4 - Suchergebnisse

424

425 Die Ergebnisse aller Beispiele entsprechen auch den Erwartungen. In allen Fällen werden
426 typischerweise als beste Ergebnisse zuerst Autos der gleichen Marke und dann Fahrzeuge mit
427 ähnlichen Merkmalen ausgewählt. Die Autos mit der geringsten Ähnlichkeit werden auch
428 angegeben (203-205). Für die Beispiele M1 und M2 werden dabei zwei komplette Autos als
429 Suchmuster verwendet. D.h. alle der Features werden im Netzwerk assoziiert und
430 anschließend mit den anderen Aktivierungsmustern verglichen. Bei den Beispielen M3-M5
431 werden nur einzelne Features selektiert und aktiviert. Im Falle von M3 wird nur Porsche
432 aktiviert. Nach Anwendung des SA Algorithmus werden aber durch das assoziative Netzwerk
433 nun auch die mit Porsche assoziierten Features aktiviert (siehe dazu das Porsche Beispiel bei
434 der Analyse Relationen). Aus diesem Grund können auch Autos anderer Hersteller gefunden
435 werden, die Werte haben, die für einen Porsche typisch sind. Bei M5 werden 3 Features
436 vorgegeben: Die Länge, die PS und MPG1. Auch hier können wieder Autos gefunden werden,
437 die am ehesten diesen Vorgaben entsprechen.

7.4 Anomalieerkennung

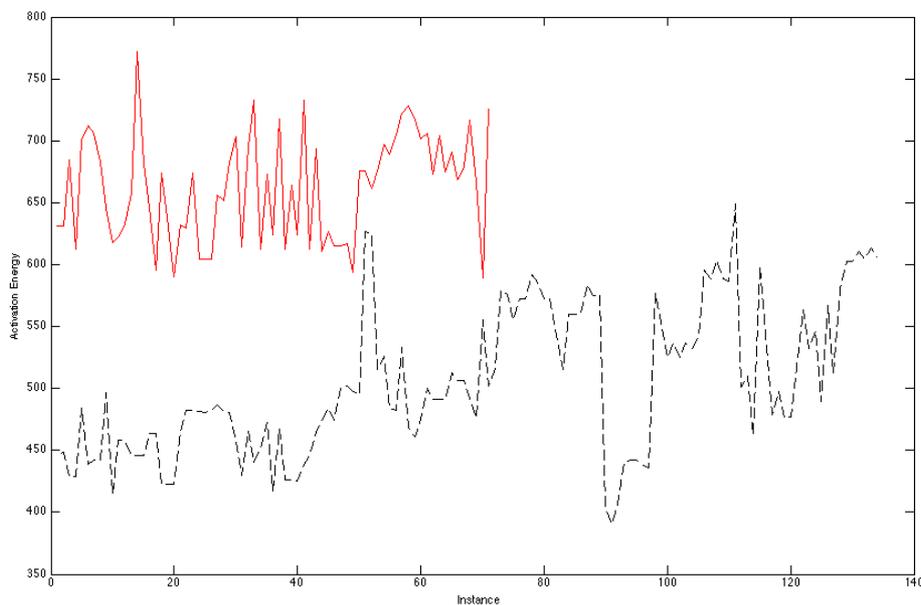
438 Hier wird gezeigt, dass das Framework auch für die Erkennung von Anomalien verwendet
439 werden kann. Im Falle von IDS macht es Sinn ein Modell für den normalen Verkehr zu

440 erstellen. Neue Muster, die während dem Training nicht vorhanden waren, und von den
 441 trainierten abweichen, werden als Anomalien bezeichnet und können z.B. Angriffe darstellen.
 442 Im Falle des hier analysierten Datensatzes wird die Anomalieerkennung anhand der zuvor
 443 gefundenen Cluster evaluiert. Dazu wird nur mit den Daten von C1 trainiert. Es handelt sich
 444 dabei um Autos unterschiedlicher Hersteller, die einen mittleren Verbrauch haben. Die
 445 anderen Cluster enthalten dann entweder stärkere Autos oder Fahrzeuge mit wenig PS und
 446 großer Reichweite. In Tabelle 5 wird die durchschnittliche Aktivierungsenergie pro Cluster
 447 (A_m) und die Standardabweichung (A_s) angegeben. Man kann hier sehen, dass die
 448 durchschnittlichen Energiewerte bei C2-C10 immer unter dem der Trainingsdaten liegen. Um
 449 dies besser darzustellen, wird in Abbildung 3 die Energie für jede Instanz der Trainingsdaten
 450 (rote Linie) und der Anomaliedaten (schwarze Linie) gezeigt. Bis auf einen kleinen Teil der
 451 Anomalien könnten diese im Prinzip durch den Vergleich mit einem festgelegten Schwellwert
 452 erkannt werden.

453

Data	A_m	A_s
Train (C1)	668	45
C2	459	40
C3	487	43
C4	473	8
C5	537	43
C6	422	20
C7	570	27
C8	524	60
C9	513	39
C10	636	10

Tabelle 5 - Aktivierungsenergie der Trainingsdaten und Anomaliedaten



454

Abbildung 3 - Aktivierungsenergie der Trainingsdaten (rot - oben) und der Anomalien (schwarz - unten)

8 Zusammenfassung und Ausblick

455 In dieser Arbeit wurde ein Framework vorgestellt, das Ereignis-Korrelation basierend auf
456 maschinellem Lernen ermöglicht. Besonderes Augenmerk wird dabei auf die unüberwachte
457 Analyse von Ereignissen gelegt. Anhand eines leicht interpretierbaren Datensatzes werden
458 die Möglichkeiten des Frameworks aufgezeigt.

459 Bei der Evaluierung des Frameworks wurden noch Punkte erkannt, die in Zukunft genauer
460 analysiert werden.

461

- 462 • **Verbesserungen beim assoziativen Netzwerk und bei den eingesetzten SA**
463 **Algorithmen:** Um die Ergebnisse des SA Algorithmus zu verbessern, müssen noch
464 weitere SA Strategien integriert werden: Beispiele dafür sind die Integration eines
465 FanOut Faktors, der den Einfluss von Knoten mit einer großen Anzahl von
466 Verbindungen zu anderen Knoten verringert. Weiters werden noch Verbesserungen
467 bei der Art der Aktivierungsfunktion benötigt.
- 468 • **Anwendung auf unterschiedliche Datensätze:** Um mehr über die Eigenschaften des
469 Frameworks zu erfahren, müssen noch weitere Datensätze aus unterschiedlichen
470 Bereichen analysiert werden.

9 Referenzen

- 471 [1] **Book**
472 Quillian, M. R. (1968). Semantic memory. In M. L. Minsky (Ed.), Semantic
473 information processing. Cambridge, MA: MIT Press
- 474 [2] **Book**
475 Fellbaum, C. (1998). Wordnet: An Electronic Lexical Database, Bradford Books
- 476 [3] **Internet reference**
477 Miller, A. et al., WordNet – a lexical database for the English language, retrieved
478 December 4, 2008 from <http://wordnet.princeton.edu/>.
- 479 [4] **Journal paper**
480 Martinetz, T.M. et al. (1991). A neural-gas network learns topologies, in T. Kohonen,
481 K. Mäkisara, O. Simula, and J. Kangas, editors, Artificial Neural Networks, pages 397-
482 402. North-Holland, Amsterdam
- 483 [5] **Book**
484 Fritzke, B. (1994). A growing neural gas network learns topologies, in Neural
485 Information Processing Systems, pages 625-632
- 486 [6] **Book**
487 Kohonen T. (1995). Self Organizing Maps, Springer
- 488 [7] **Book**
489 McLachlan G. (2008). The EM Algorithm and Extensions, 2nd Edition, Wiley
- 490 [8] **Chapter in proceedings**
491 Hamp, B. (1997). GermaNet – a Lexical-Semantic Net for German, in Proceedings of
492 ACL workshop Automatic Information Extaction and Building of Lexical Semantic
493 Resources for NLP Applications, Madrid
- 494 [9] **Website**
495 Asuncion, A., Newman, D. (2007). UCI machine learning repository, retrieved
496 February 6, 2009 from <http://archive.ics.uci.edu/ml/>.

10 Anhänge

497 Alle Ergebnisse sind in den folgenden beigelegten Dateien notiert. Bei den relevanten Stellen im
498 Dokument wird auf diese Dokumente referenziert.

499 ***test-a.csv***

500 ***test-b.csv***

501 ***test-c.csv***

502 ***test-d.csv***

503 ***test-e.csv***

504 ***test-f.csv***

505 ***test-g.csv***

506 ***cluster.csv***

Historie

Version	Datum	Kommentar
1.0	06.02.2009	Finale Version
Ersteller Peter Teufl		

507