

# Implementierung datenbasierter Verfahren zur Prozessüberwachung und Fehlerdiagnose in ein Softwaresystem für das Anlagenmanagement

Gunar Gruss, Hochschule Lausitz (FH)

gunar.gruss@hs-lausitz.de

Peter Engel, PC-Soft GmbH

pengel@pcsoft.de

Torsten Jeinsch, Hochschule Lausitz (FH)

torsten.jeinsch@hs-lausitz.de

Steven X. Ding, Universität Duisburg-Essen

steven.ding@uni-due.de

## Zusammenfassung

Datengestützte Verfahren gewinnen in der Überwachung von Anlagen immer mehr an Bedeutung. Ihr Vorteil besteht darin, dass direkt aus den Prozessmessdaten ein Modell entwickelt wird ohne Kenntnis spezifischer regelungstechnischer Zusammenhänge. Im ZIM-Projekt *Entwicklung eines neuen adaptiven Verfahrens und Systems für das technische Anlagenmanagement von Industrieanlagen (EVA)* werden datenbasierte Verfahren in das Anlagenmanagementsystem zedas<sup>®</sup><sup>1</sup> integriert. Die effektive Validierung neuer Parametrierungen der Verfahren in der Praxis erfolgt mit Hilfe eines Engineeringmoduls. Diese Vorgehensweise wird exemplarisch für die Hauptkomponentenanalyse beschrieben.

## 1 Einleitung

In der Industrie werden immer komplexere, hoch automatisierte Maschinen und Anlagen eingesetzt. Dies hat zur Folge, dass die physikalisch-technische Interpretation von Messwerten der Anlagenparameter immer schwieriger wird. Um fehlerhafte Zustände frühzeitig erkennen zu können, verlangt die Industrie nach Werkzeugen die Messwerte aufzubereiten, zu interpretieren und in geeigneter Weise darzustellen. Für die Bereitstellung dieser Werkzeuge werden die im BMBF-Projekt *Intelligente Verfahren zur Diagnose und Prozessführung* untersuchten multivariaten statistischen Verfahren, wie die Hauptkomponentenanalyse (engl. Principal Component Analysis - PCA) in das Softwaresystem zedas<sup>®</sup>

---

<sup>1</sup>Marke der PC-Soft GmbH

für das Anlagenmanagement von komplexen automatisierungstechnischen Systemen integriert. Die Evaluierung erfolgt an automatisierungstechnischen Anlagen in der Praxis.

## 2 Datenbasierte Verfahren zur Prozessüberwachung und Fehlerdiagnose

Methoden der schließenden Statistik erlauben es aus einer Stichprobe mit allgemeinen Gesetzmäßigkeiten auf Eigenschaften einer Grundgesamtheit zu schließen. [9] Den Ausgangspunkt datenbasierter Verfahren bilden historische Prozessmessdaten. Für weitere Betrachtungen gehen wir davon aus, dass diese in einer Matrix  $X_0, X_0 \in R^{n \times m}$ , organisiert sind, wobei  $m$  die Anzahl der Datenspuren und  $n$  die Anzahl der Samples abbildet. Für die Anwendung multivariater Verfahren in der Praxis ist eine Vorverarbeitung der Daten notwendig. Je nach Konzept der Datenspeicherung werden erfasste Messwerte nur abgespeichert, wenn sich der Betrag oder das Verhältnis zum Vorgänger um einen bestimmten Betrag ändert. Fehler in der Datenübertragung können auch zu unplausiblen Einträgen in den Datenbanken führen. Aus diesen Bedingungen ergeben sich folgende Schritte der Datenvorverarbeitung:

1. Datenvektoren mit keinen oder zu wenig Einträgen werden entfernt
2. Interpolation bzw. Extrapolation in Bereichen mit fehlenden Daten
3. Beseitigung von Ausreißern

Aus Prozessmessdaten, welche den fehlerfreien Prozesszustand charakterisieren wird ein statistisches Prozessmodell erstellt. Zu diesem gelangt man durch statistische Methoden mit denen Parameter des Modells ermittelt werden. Dieser Teil der Prozessüberwachung und Fehlerdiagnose wird als Modellbildung oder Trainingsphase bezeichnet und erfolgt in der Regel offline. Die Überwachung wird bevorzugt online durchgeführt. Dabei

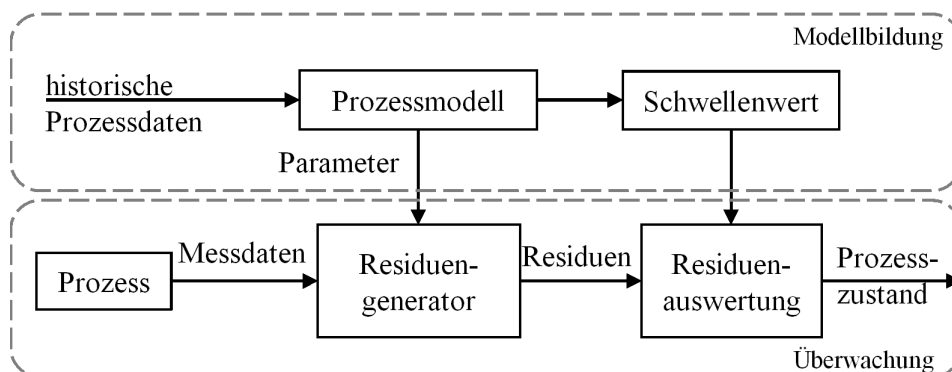


Abbildung 1: Grundprinzip der Fehlerdiagnosemethoden

werden die Eigenschaften der aktuellen Prozessmessdaten (Istdaten) mit den Modellparametern verknüpft. Dieser Vorgang wird als Residuengenerierung bezeichnet. Im nächsten

Schritt werden die Residueninformationen mit den in der Prozessmodellbildung ermittelten Schwellwerten verglichen. Neben der Fehlerdetektion kann auch die Fehlerisolation Bestandteil der Residuenauswertung sein. In Abbildung 1 ist das Prinzip der Fehlerdiagnosemethoden dargestellt [4][5].

Ein häufig eingesetztes Verfahren stellt die Hauptkomponentenanalyse dar.

## 2.1 Hauptkomponentenanalyse

Nachfolgend werden wesentliche Grundlagen zur Theorie der Hauptkomponentenanalyse erläutert. Primäres Ziel ist die Entdeckung von Zusammenhängen zwischen Variablen, Objekten oder Merkmalen. Zu Beginn der Analyse sind die Beziehungen der Variablen des Datensatzes untereinander unbekannt. Mit Hilfe der Hauptkomponentenanalyse werden aus den Variablen wenige latente, voneinander unabhängige Faktoren extrahiert. Die großen Vorteile dieses Verfahrens liegen darin, dass es

- analytisch fassbar ist,
- mit einer großen Anzahl von Prozessvariablen funktioniert und
- relativ geringen Rechenaufwand verursacht.

Den Ausgangspunkt datenbasierter Verfahren bilden historische Prozessmessdaten. Aus der normalisierten Datenmatrix  $X$ ,

$$X = \frac{X_0 - \mu}{\sigma}, \quad (1)$$

mit dem Vektor der Mittelwerte der Spalten  $\mu$  und der Standardabweichung der Spalten  $\sigma$ , wird die Varianz-Kovarianzmatrix  $\Sigma$  gebildet:

$$\Sigma = \frac{1}{n-1} X^T X. \quad (2)$$

Mit der Singular Value Decomposition (SVD) werden die Eigenvektoren und Eigenwerte bestimmt. Als Ergebnis erhält man die Diagonalmatrix  $\Lambda$  mit den nicht negativen realen Eigenwerten in absteigender Ordnung,  $\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)$ , sowie die Matrix  $P \in R^{m \times m}$  mit den zugehörigen Eigenvektoren.

$$\frac{1}{n-1} X^T X = P \Lambda P^T \quad (3)$$

Die Eigenvektoren mit den größten Eigenwerten spiegeln den größten Teil der Variationen wieder. Sie bezeichnet man als Hauptkomponenten. Der  $m \times m$ -dimensionale Raum  $P$  kann in einen Unterraum der Hauptkomponenten und in einen Residuum-Unterraum aufgeteilt werden.

$$P = [P_{pc} P_{res}] \in R^{m \times m}, P_{pc} \in R^{m \times l}, P_{res} \in R^{(m-l) \times (m-l)} \quad (4)$$

Mit Hilfe der Loading-Matrix  $P$  werden die Originaldaten in den Unterraum der Hauptkomponenten projiziert. Als Ergebnis erhält man die Score-Matrix  $T$ .

$$T = XP \quad (5)$$

Zur Bestimmung der Dimension des Raumes der Hauptkomponenten gibt es verschiedene Ansätze. Die zwei bekanntesten sind [6]:

- Kriterium von Kaiser

Alle Eigenwerte größer als der Mittelwert aller Eigenwerte bilden die Hauptkomponenten (1).

$$l = \max \{i | \lambda_i^2 \geq \bar{\lambda}^2\} \quad (6)$$

- Kriterium von Jolliffe

Jolliffe hat festgestellt, dass die Forderung von Kaiser zu hoch ist und bestimmt, dass alle Eigenwerte die größer als das 0,7-fache des Mittelwerts sind die Größe des Raumes der Hauptkomponenten bestimmen.

$$l = \min \left\{ m \mid \lambda_1^2 + \dots + \lambda_m^2 \geq 0,7 \lambda_{gesamt}^2 \right\} \quad (7)$$

## 2.2 Fehlererkennung

Für die Erkennung von Fehlern mit Hilfe multivariater Verfahren und insbesondere der Hauptkomponentenanalyse werden im wesentlichen zwei Kennwerte [2] genutzt:

- Hotelings  $T^2$ -Statistik,
- $SPE$ -Statistik (engl. Squared Prediction Error).

### $T^2$ -Statistik

Die  $T^2$ -Statistik folgt der  $F$ -Verteilung und ist von zwei Freiheitsgraden abhängig, der Anzahl der Trainingssample ( $n$ ) und der Anzahl der Hauptkomponenten ( $l$ ).

$$T_\alpha^2 = \frac{l(n-1)(n+1)}{n(n-l)} F_\alpha(l, n-l), \quad (8)$$

wobei  $\alpha$  das Signifikanzniveau angibt.

Um den aktuellen Samplevektor zu bewerten, wird dieser mit Hilfe der Eigenvektoren und Eigenwerte des Hauptkomponentenraumes in diesen projiziert. Man erhält den  $T^2$ -Index.

$$T^2 = x^T P \Lambda P^T x \quad (9)$$

## *SPE*- oder *Q*-Statistik

Der Grenzwert für die *SPE*-Statistik wurde von Jackson und Mudholkar [8] approximiert:

$$Q_\alpha = \Theta_1 \left( \frac{c_\alpha \sqrt{2\Theta_2 h_0^2}}{\Theta_1} + 1 + \frac{\Theta_2 h_0 (h_0 - 1)}{\Theta_1^2} \right)^{\frac{1}{h_0}} \quad (10)$$

mit

$$\Theta_i = \sum_{j=l+1}^m \lambda_j^i, i = 1, 2, 3 \quad (11)$$

und

$$h_0 = 1 - \frac{2\Theta_1\Theta_3}{3\Theta_2^2}, \quad (12)$$

wobei  $j$  dem Index der Hauptkomponenten und  $c_\alpha$  dem Wert der Normalverteilung zum Signifikanzniveau  $\alpha$  entspricht.

Zur Bewertung der Eigenschaften des aktuellen Samples im Residuenraum wird der *SPE*-Indexwert berechnet.

$$Q = SPE = \|(I - PP^T)\|^2 x = x^T (I - PP^T)^2 x \quad (13)$$

Ist für einen Samplevektor einer oder beide der berechneten Indices größer als der Grenzwert, entspricht der Zustand der Anlage zum Konfidenzintervall  $\alpha$  nicht dem Normalzustand.

## 3 Implementierung datenbasierter Verfahren

Im Zuge der Weiterentwicklung des Softwaresystems für das Anlagenmanagement werden datenbasierte Verfahren für die Prozesszustanderkennung und Fehleridentifizierung integriert.

Um neue Verfahren in der Praxis testen und später eine kontinuierliche Weiterentwicklung der Verfahren durchführen zu können, werden diese in einem Engineeringmodul zusammengefasst, welches im Hauptsystem eingebettet ist. Im Gegensatz zu diesem ist das Engineeringmodul in der Skriptsprache *R* realisiert. Das Open-Source-Projekt *R* basiert auf *S*, der Urversion der interaktiven Softwarelösung für statistische Anwendungen *S – PLUS* [1]. Diese Vorgehensweise bietet folgende Vorteile:

- ungefilterter Zugriff auf Produktionsdaten,
- schnelle Kommunikation mit den für den Prozess verantwortlichen Personen bei Fehlerdetektion bzw. Missdetektion

- schneller Lernzyklus für die Parametrierung der entwickelten Algorithmen,
- erkennen von Anforderungen welche die Praxis stellt,
- Eigenständigkeit, aber auch Nutzung der Ressourcen des Hauptsystems,

verbunden mit Anforderungen nach:

- Robustheit gegenüber Fehlfunktion → das Hauptsystem darf nicht beeinträchtigt werden,
- effektiver Datenvorverarbeitung,
- Überwachung der Ergebnisse und ggf. Neuparametrierung der Algorithmen und
- zeitnaher Implementierung neuer Parametrierungen.

Durch den Einsatz des Engineeringmoduls werden Algorithmen parametriert oder neu integriert, ohne dass die Beherrschung einer Hochsprache noch das Kompilieren neuen Softwarecodes unter den jeweils spezifischen Randbedingungen des Systems erforderlich ist, d.h. der Entwickler kann sich voll auf die Entwicklung der Algorithmen der Verfahren konzentrieren. Um unbefugten Zugriff auf das Engineeringmodul, z.B. durch einen Operator zu verhindern, ist es durch ein Passwort geschützt.

Der Zugriff auf Daten erfolgt unter Nutzung von Funktionen des Hauptsystems. Somit

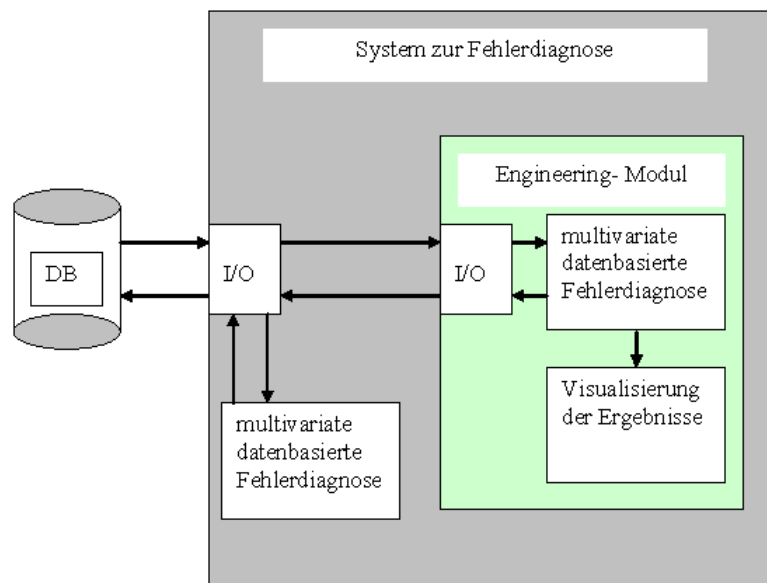


Abbildung 2: Einbettung des Engineeringmoduls

können alle Daten, welche für die Standardmodule verfügbar sind, auch für das Engineeringmodul genutzt werden. Das Engineeringmodul gliedert sich in die Bereiche Datenvorverarbeitung, Modellbildung (Training), Überwachung.

### 3.1 Datenvorverarbeitung

Nach der Auswahl und dem Einlesen der Datenspuren erfolgt die Vorverarbeitung der Daten. Die Datenvektoren werden zu einer Matrix zusammengefasst, wobei leere oder nur mit wenigen Werten besetzte Datenvektoren nicht berücksichtigt werden. Für nicht besetzte Stellen werden durch Inter- bzw. Extrapolation Werte erzeugt. Für die nun vollbestzte Datenmatrix erfolgt die Ermittlung und Beseitigung von Ausreißern. Die vorverarbeitete Datenmatrix ist die Grundlage für die Modellbildung.

### 3.2 Modellbildung

Die Aufgabe der Modellbildung ist die Ermittlung statistischer Kennwerte. Anhand dieser Kennwerte erfolgt später die Überwachung. Für verschiedene Verfahren der datenbasierten Fehlerdiagnose unterscheidet sich die Modellbildung. Für das Beispiel Hauptkomponentenanalyse ist der Ablauf in Abbildung 3a) dargestellt.

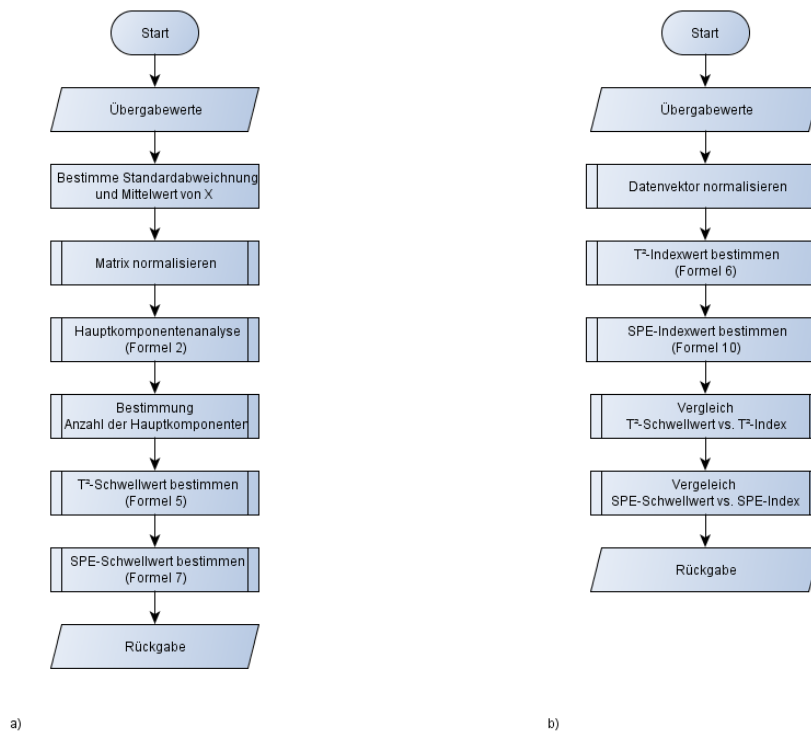


Abbildung 3: a) PCA-Modellbildung b) PCA-Überwachung

### 3.3 Überwachung

Typischerweise erfolgt die Überwachung online, kann aber für Testzwecke auch offline durchgeführt werden. Nach dem Einlesen und Normalisieren der aktuellen Prozessdaten werden, wie im Punkt 2.2 beschrieben, für jeden Abtastzeitpunkt die  $T^2$ - und  $SPE$ -Indices berechnet. Überschreitet mindestens einer der beiden berechneten Indices den

jeweiligen Schwellwert wird ein Alarm generiert.

$$Alarm = T^2 > T_\alpha^2 \vee SPE > Q_\alpha \quad (14)$$

Den Ablauf der Überwachung verdeutlicht Abbildung 3b). Für ein Gurtfördersystem,

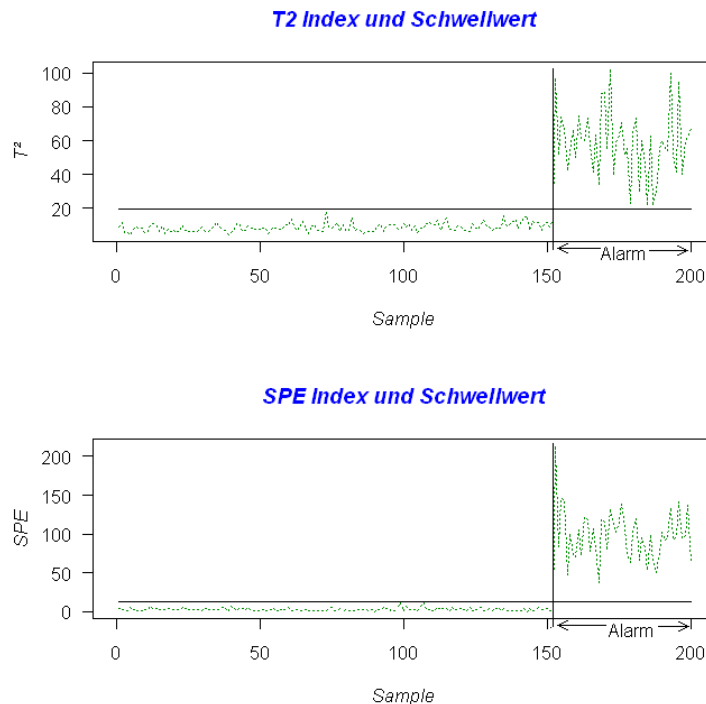


Abbildung 4:  $T^2$ - und  $SPE$ -Index

welches als Pilotanwendung dient, ist in Abbildung 4 ein Ausschnitt des Verlaufs beider Indices dargestellt. Deutlich zu erkennen ist die Überschreitung beider Schwellwerte.

### 3.4 Auswahl geeigneter Verfahren

Zum Vergleich der Leistungsfähigkeit der datenbasierten Verfahren im Standardmodul und Engineeringmodul werden ermittelte Kennwerte (z.B. Grenzwerte und Indices) aus der Datenbank ausgelesen und bewertet. Als Kriterien für die Bewertung werden die Detektions- bzw. Missdetektionsrate herangezogen. Die Detektionsrate gibt an, wie viele Ereignisse, welche nicht dem Normalzustand entsprechen, durch den Algorithmus gefunden werden. Im Gegensatz dazu sagt die Missdetektionsrate aus, wie viele Ereignisse als Fehler erkannt werden, obwohl der Normalzustand vorliegt. Nach erfolgter Erprobung neuartiger bzw. neu parametrierter Verfahren werden diese in das Hauptmodul integriert und stehen dann dem Anlagenoperator zur Verfügung.

## Literatur

- [1] R Development Core Team: *R: A Language and Environment for Statistical Computing*, <http://www.R-project.org>, R Foundation for Statistical Computing, 2011.



- [2] Russell, Evan L.; Chiang, Leo H.; Braatz, Richard D.: *Data driven methods for fault detection and diagnosis in chemical processes*, Springer-Verlag, 2000.
- [3] Sachs, Lothar; Heddrich, Jürgen: *Angewandte Statistik - Methodensammlung mit R*, Springer-Verlag, 2006.
- [4] Ding, Steven X.: *Model-based Fault Diagnosis Techniques* Springer-Verlag 2008
- [5] Lunze, Jan: *Automatisierungstechnik: Methoden für die Überwachung und Steuerung kontinuierlicher und ereignisdiskreter Systeme* Oldenburg Wissenschaftsverlag GmbH 2008
- [6] Handl, Andreas: *Multivariate Verfahren Theorie und Praxis unter besonderer Berücksichtigung von S-Plus*, Springer-Verlag, 2002.
- [7] Engel, Peter; Jeinsch, Torsten; Schoch, Daniel; Stargala, Torsten; Ding, Steven X.: *Entwicklung selbstlernender datenbasierter Verfahren zur Überwachung und Diagnose von Industrieanlagen*, VDI-Berichte 2092, VDI Verlag, 2010.
- [8] Jackson, Edward J.; Mudholkar, Govind S.: *Control Procedures for Residuals Associated with Principal Component Analysis*, *Technometrics*, Vol. 21, No. 3 (Aug., 1979), pp. 341-349.
- [9] Fahrmeir, Ludwig; Hamerle Alfred, Tutz Gerhard: *Multivariate statistische Verfahren*, Walter de Gruyter, 1996