

Réseaux bayésiens pour l'identification de variables hors-contrôle

Sylvain Verron, Teodor Tiplica, Abdessamad Kobi

► **To cite this version:**

Sylvain Verron, Teodor Tiplica, Abdessamad Kobi. Réseaux bayésiens pour l'identification de variables hors-contrôle. 5ème Conférence Internationale Francophone d'Automatique (CIFA'08), 2008, Bucarest, Roumanie. inria-00517046

HAL Id: inria-00517046

<https://hal.inria.fr/inria-00517046>

Submitted on 13 Sep 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Réseaux bayésiens pour l'identification de variables hors-contrôle

Sylvain VERRON, Teodor TIPLICA, Abdessamad KOBİ

Laboratoire en Sûreté de Fonctionnement, Qualité et Organisation (LASQUO)
ISTIA, 62 Avenue Notre Dame du Lac, 49000 ANGERS, France

sylvain.verron@univ-angers.fr

<http://www.univ-angers.fr/laboratoire.asp?ID=34&langue=1>

Résumé—Le but de cet article est de présenter une méthode de détection et d'identification par réseaux bayésiens. Pour cela, une combinaison est réalisée entre les récents travaux de Li et al. [1] (décomposition causale du T^2) et certains de nos précédents travaux [2], [3] (cartes de contrôle multivariées par réseaux bayésiens). Ainsi, pour un procédé multivarié, les améliorations proposées permettent à la fois la détection d'une faute et l'identification des variables impliquées dans celle-ci. Un intérêt particulier de cette méthode réside dans le fait qu'elle n'exploite qu'un seul et même outil : un réseau bayésien.

Mots-clés—Réseaux bayésiens, Cartes de contrôle, Décomposition MYT.

I. INTRODUCTION

De nos jours, les procédés industriels possèdent de plus en plus de capteurs, fournissant ainsi une importante quantité de données. Un champ de recherche intéressant porte sur l'utilisation de ces données pour contrôler le procédé. Le contrôle d'un procédé peut être vu comme une procédure en 4 phases [4]. Dans la première phase, la détection, l'objectif est de détecter une situation anormale, une faute dans le procédé. Le but de la seconde phase, l'identification de faute, est d'identifier les variables les plus significatives pour le diagnostic de la faute. La troisième phase est le diagnostic de faute, elle consiste à déterminer quel type de faute est apparue dans le procédé. Finalement, la dernière phase est celle de la reconfiguration du procédé qui permet d'agir sur le procédé ou sa commande afin de retrouver les conditions nominales d'utilisation.

Le contrôle des procédés peut être réalisé par trois principales approches [4] : l'approche analytique, l'approche basée sur les données, et l'approche à base de connaissances. L'approche analytique consiste à construire un modèle mathématique du procédé. L'approche à base de connaissances se base sur des modèles qualitatifs du procédé. Enfin, les méthodes basées sur les données exploitent des développements statistiques des données du procédé. Théoriquement, les méthodes analytiques sont les plus rigoureuses et les plus précises, elles sont donc celles donnant les meilleurs résultats. Cependant, pour des systèmes importants et complexes (nombre élevé d'entrées, de sorties, d'états de fonctionnement), l'obtention de modèles assez détaillés est extrêmement difficile. Ainsi, les méthodes analytiques ne sont pas les mieux adaptées pour ce type de système, pouvant mener à des conclusions erronées. Dans ce cas, les méthodes basées sur des

développements statistiques rigoureux seront préférées aux méthodes analytiques.

La littérature est riche concernant les techniques de contrôle basées sur les données : maîtrise statistique des procédés univariés (cartes de contrôle de Shewhart) [5], [6], maîtrise statistique des procédés multivariés (cartes de contrôle T^2 et Q) [7], [8], ainsi que des techniques basées sur l'Analyse en Composantes Principales (ACP) [9] telles que l'ACP multiéchelle ou l'ACP dynamique [10]. Kano et al. [11] effectuent une comparaison de ces différentes méthodes. D'autres méthodes se basent sur la Projection dans les Structures Latentes (PSL) [12] comme la PSL multiéchelle [13]. Concernant l'identification de fautes, Tiplica et al. [14] ont établi un comparatif de plusieurs méthodes. L'une des techniques statistiques la plus intéressante est la décomposition MYT [15], [16] qui effectue une décomposition de la statistique T^2 en composantes orthogonales permettant de déterminer quelle variable ou groupe de variables a contribué à une situation hors-contrôle (faute). Récemment, Li et al. [1] ont proposé une amélioration de la décomposition MYT nommée décomposition causale du T^2 . Pour cela, ces auteurs se basent sur un réseau bayésien causal représentant les différentes variables du procédé.

Comme nous l'avons présenté, la détection et l'identification de fautes se basent sur des différents outils (carte de contrôle, décompositions diverses, etc). Il serait intéressant d'essayer de regrouper tous ces méthodes sous un seul et même outil. L'objectif de cet article est de proposer une amélioration de la méthode proposée par Li et al. [1], afin de n'utiliser qu'un seul et même réseau bayésien, que ce soit pour la détection de faute ou pour l'identification des variables impliquées dans cette faute. Suite à une présentation des réseaux bayésiens (section II), la section III permet de détailler les méthodes de décomposition MYT et causale. La section IV rappelle tout d'abord les principes de construction d'une carte de contrôle multivariée par réseaux bayésiens, puis présente les moyens d'obtenir la détection et l'identification de situations hors-contrôle dans un seul et même réseau bayésien. Une application de nos propositions sur un exemple simple est présenté à la section V, puis les conclusions et perspectives seront énoncés dans la section VI.

II. RÉSEAUX BAYÉSIENS

Un Réseau Bayésien (RB) [17], [18] est un modèle graphique dans lequel les connaissances sont représentées sous forme de variable. Chaque variable est un nœud du graphe et prend ses valeurs dans un ensemble discret ou continu. Le graphe est toujours dirigé et acyclique. Les arcs dirigés représentent un lien de dépendance directe (la plupart du temps il s'agit de causalité). Ainsi un arc allant de la variable X à la variable Y exprimera le fait que X dépend directement de Y . L'absence d'arc ne renseigne alors que sur la non-existence d'une dépendance directe. Les paramètres expriment les poids donnés à ces relations et sont les probabilités conditionnelles des variables sachant leurs parents (exemple : $P(Y|X)$) ou les probabilités a priori si la variable n'a pas de parents. On peut définir formellement un réseau bayésien comme étant un triplet $\{\mathbf{G}, \mathbf{E}, \mathbf{D}\}$ où :

$\{\mathbf{G}\}$ est un graphe acyclique orienté, $\mathbf{G} = (V, A)$, où V est l'ensemble des nœuds de \mathbf{G} , et A est l'ensemble des arcs de \mathbf{G} ,

$\{\mathbf{E}\}$ est un espace probabilisé fini (Ω, Z, P) , avec Ω un espace non vide, Z un ensemble de sous-espaces de Ω , et P une mesure de probabilité sur Z avec $P(\Omega) = 1$,

$\{\mathbf{D}\}$ est un ensemble de variables aléatoires associées aux nœuds de \mathbf{G} et défini sur \mathbf{E} , tel que :

$$P(V_1, V_2, \dots, V_n) = \prod_{i=1}^n P(V_i | C(V_i)) \quad (1)$$

où $C(V_i)$ est l'ensemble des causes (parents) de V_i dans le graphe \mathbf{G} .

Il est possible de réaliser des classifieurs performants grâce aux réseaux bayésiens [19], [20], [21]. Nous présentons ici les principaux types de structures pour employer les réseaux bayésiens comme classifieurs. Un classifieur bayésien d'un problème à p variables a pour particularité de posséder $p+1$ nœuds. En effet, tous les classifieurs bayésiens modélisent le fait d'appartenance à une classe par un nœud discret. Nous nommons ce nœud "nœud de classe", et nous le notons C . Ce nœud est un nœud discret multinomial à k modalités, où k représente le nombre de classes de notre problème (C_1, C_2, \dots, C_k) . Ce nœud de classe a pour particularité de ne pas posséder de nœud parent. Les autres variables, au nombre de p , que nous nommons variables descriptives, sont notées X_i (i de 1 à p).

Le classifieur bayésien possédant la structure la plus simple est le Réseau Bayésien Naïf (RBN), également appelé classifieur de Bayes (figure 1 (a)). On le qualifie de naïf car il fait l'hypothèse, très forte, que chaque variable descriptive est, conditionnellement à la classe, indépendante des autres. Lorsque toutes les variables descriptives sont incorporées au modèle, on parlera alors de structure naïve complète. Ce classifieur est extrêmement connu car ses performances (notamment dans le cas où toutes les variables sont discrètes) sont intéressantes dans certains domaines et dépassent des techniques beaucoup plus sophistiquées même lorsque l'hypothèse d'indépendance est violée [22]. Au vu de l'hypothèse forte que ce classifieur implique, il est normal que certains chercheurs aient voulu améliorer ce classifieur. Friedman [20] propose d'ajouter des arcs entre

les différentes variables descriptives du classifieur naïf. Pour cela, il décide de créer un arbre entre les variables descriptives, à la manière de Chow et Liu [23], afin d'obtenir un TAN (Tree Augmented Naïve Bayes), visible sur la figure 1 (b). L'algorithme part d'un réseau bayésien naïf et ajoute alors un arc entre les variables qui partagent la plus importante information mutuelle. Mais, pour respecter la topologie de l'arbre, l'algorithme interdit à chaque nœud d'avoir plus de 2 parents (soit un parent en plus du nœud classe). Afin de prendre en compte la corrélation entre les différents descripteurs, il a également été proposé le CSNBN (Condensed Semi-Naïve Bayesian Network) : un réseau bayésien semi naïf condensé [24], [25] (figure 1 (c)). On le nomme condensé car il introduit un nouveau type de variable : les variables jointes. Ces nouvelles variables jointes représentent en fait un groupement de variables descriptives. Bien entendu, une variable descriptive ne pourra se trouver que dans une seule variable jointe. Le fait que deux variables se trouvent dans une variable jointe implique que ces deux variables sont corrélées. Un regroupement de p variables continues suivra une loi normale multivariée et sera donc représenté par un seul nœud continu de dimension p .

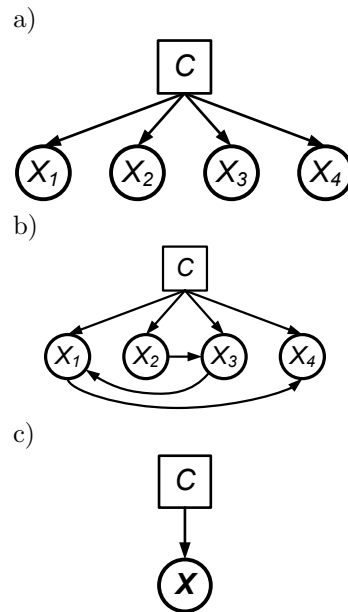


Fig. 1. Différents classifieurs bayésiens : RBN (a), TAN (b) et CSNBN (c).

III. DIFFÉRENTES MÉTHODES

A. Diagnostic par décomposition MYT

Comme nous l'avons déjà dit, un principe de détection pour les procédés multivariés est la carte de contrôle du T^2 , mise au point par Hotelling [7]. Cependant, la carte ne donne aucune indication concernant le diagnostic de la situation hors-contrôle. Pour cela, beaucoup de méthodes ont été proposées [26], [27], [28], [15], [29]. Une étude comparative est effectuée par Tiplica [30].

Une décomposition du T^2 particulièrement intéressante a été mise au point par Mason, Young et Tracy [15], d'où le nom "décomposition MYT". De plus, pour com-

prendre cette méthode de manière plus intuitive, les auteurs donnent un exemple avec un procédé bivarié [16]. Il est également à préciser que les auteurs ont prouvé que certaines méthodes peuvent se ramener à des cas particuliers de décomposition MYT [15]. En effet, cette décomposition réunit les idées de Hawkins [28], basées sur la régression multiple et l'analyse des résidus, et de Doganaksoy [27] sur la contribution des variables à la statistique de Student.

Le principe de la méthode MYT est de décomposer la statistique T^2 dans un nombre limité de composantes orthogonales qui sont également des distances statistiques (et donc surveillables). La décomposition est la suivante :

$$T^2 = T_1^2 + T_{2\bullet 1}^2 + T_{3\bullet 1,2}^2 + T_{4\bullet 1,2,3}^2 + \dots + T_{p\bullet 1,2,3\dots p-1}^2 \quad (2)$$

où $T_{i\bullet j,k}^2$ représente la statistique T^2 de la régression des variables X_j et X_k sur la variable X_i . On voit qu'il existe un nombre important de décompositions différentes ($p!$), et donc qu'il existe un grand nombre de facteur ($p \times 2^{p-1}$) différents. Pour mieux comprendre, sur un procédé à 3 variables, nous obtenons les différentes décompositions suivantes :

$$\begin{aligned} T^2 &= T_1^2 + T_{2\bullet 1}^2 + T_{3\bullet 1,2}^2 \\ T^2 &= T_1^2 + T_{3\bullet 1}^2 + T_{2\bullet 1,3}^2 \\ T^2 &= T_2^2 + T_{1\bullet 2}^2 + T_{3\bullet 1,2}^2 \\ T^2 &= T_2^2 + T_{3\bullet 2}^2 + T_{1\bullet 2,3}^2 \\ T^2 &= T_3^2 + T_{1\bullet 3}^2 + T_{2\bullet 1,3}^2 \\ T^2 &= T_3^2 + T_{2\bullet 3}^2 + T_{1\bullet 2,3}^2 \end{aligned} \quad (3)$$

Le calcul des termes n'est pas détaillé ici, mais on pourra bien entendu se reporter aux travaux de Mason et al. [15], [16]. Il est à noter que les termes T_j^2 sont appelés facteurs non-conditionnés (puisque'ils ne dépendent pas du tout des autres variables que j), alors que les autres termes sont appelés facteurs conditionnés. Ce qui est intéressant, c'est que chaque facteur suit une distribution de Fisher (à une constante multiplicative près) :

$$T_{j+1\bullet 1,\dots,j}^2 = \frac{(m+1)(m-1)}{m(m-k-1)} F_{1,m-k-1} \quad (4)$$

où k est le nombre de facteurs conditionnés. On pourra donc simplifier cette équation pour les termes non-conditionnés ($k=0$) par :

$$T_{j+1\bullet 1,\dots,j}^2 \sim \frac{m+1}{m} F_{1,m-1} \quad (5)$$

Cela nous permet de détecter un problème sur chacun des facteurs de la décomposition. Par exemple, si l'on s'aperçoit que le facteur $T_{2\bullet 1}^2$ est responsable d'un hors contrôle du procédé, on peut immédiatement aller chercher la cause de l'anomalie sur un réglage physique affectant la corrélation entre ces deux variables. Mais, pour moins de calcul, il suffit d'utiliser une carte T^2 pour la détection de situation hors-contrôle, et si une faute se produit, alors on utilise la méthode MYT pour déterminer d'où vient cette faute. L'analyse des facteurs se fait dans l'ordre de niveau (ex : T_1^2, T_2^2, T_3^2 , puis $T_{2\bullet 1}^2, T_{3\bullet 1}^2, T_{1\bullet 2}^2, T_{3\bullet 2}^2, T_{1\bullet 3}^2, T_{2\bullet 3}^2$ puis finalement $T_{3\bullet 1,2}^2, T_{2\bullet 1,3}^2, T_{1\bullet 2,3}^2$) jusqu'à ce qu'on trouve le facteur ayant causé la détection d'une erreur sur la carte T^2 .

L'avantage de la méthode MYT est qu'elle fournit un diagnostic d'une situation hors-contrôle, sans avoir à la comparer à des exemples de fautes préalablement apparues dans le procédé. Ainsi, cette méthode de diagnostic est une méthode non-supervisée. De plus, un autre avantage de cette méthode est qu'elle est basée sur la même démarche que la carte T^2 . Les outils statistiques sont les mêmes et on peut penser qu'une implémentation pratique est beaucoup plus compréhensible qu'un mélange de plusieurs techniques.

B. Décomposition causale du T^2

La méthode MYT est très intéressante mais elle est sujette à un inconvénient majeur : le nombre de termes à calculer. En effet, elle impose un nombre de décompositions égale à $p!$ (où p est le nombre de variables du procédé). Or, ces décompositions impliquent alors le calcul de $p \times 2^{p-1}$ termes distincts. Par exemple, pour un procédé à 20 variables, plus de 10 millions de termes distincts sont à calculer. Des efforts furent effectués afin de réduire le nombre de termes en appliquant un algorithme en 5 étapes [16] permettant une réduction significative du nombre de termes à calculer. Cependant, Li et al. [1] font la remarque que même avec l'utilisation de l'algorithme, le nombre de termes à calculer est tout de même important (très supérieur à p , notamment en présence de fautes multiples). Li et al. [1] proposent alors une méthode exploitant les réseaux bayésiens : la décomposition causale du T^2 . Un graphe causal représentant le procédé permet de réduire le nombre de termes à calculer à p . En plus de la diminution de calcul engendrée par cette méthode, les auteurs précisent que les performances sont également améliorées.

L'hypothèse de base de la méthode proposée par Li et al. [1] est que le procédé peut être modélisé sous la forme d'un réseau bayésien causal où chaque variable du procédé est une variable gaussienne univariée. Lorsqu'un réseau bayésien représente uniquement des variables continues normales, il est également appelé modèle linéaire gaussien. Ainsi, pour un procédé à 3 variables, on peut par exemple obtenir le réseau bayésien de la figure 2.

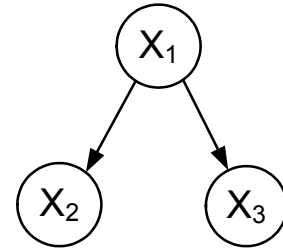


Fig. 2. Exemple d'un modèle causal linéaire gaussien

Dans le cadre de la modélisation du procédé par un modèle linéaire gaussien, les auteurs font la distinction entre deux types de décomposition MYT : "pour une décomposition du T^2 donnée, s'il existe un terme $T_{i\bullet 1,\dots,i-1}^2$ tel que l'ensemble de variables $\{X_1, \dots, X_{i-1}\}$ contient au moins un descendant de X_i , alors cette décomposition est de type A, dans le cas contraire, la décomposition est de type B". Ainsi, nous pouvons classer dans la table I les différentes décompositions du procédé à 3 variables de la figure 2.

Décomposition	Type
$T^2 = T_1^2 + T_{2\bullet 1}^2 + T_{3\bullet 1,2}^2$	Type B
$T^2 = T_1^2 + T_{3\bullet 1}^2 + T_{2\bullet 1,3}^2$	Type B
$T^2 = T_2^2 + T_{1\bullet 2}^2 + T_{3\bullet 1,2}^2$	Type A
$T^2 = T_2^2 + T_{3\bullet 2}^2 + T_{1\bullet 2,3}^2$	Type A
$T^2 = T_3^2 + T_{1\bullet 3}^2 + T_{2\bullet 1,3}^2$	Type A
$T^2 = T_3^2 + T_{2\bullet 3}^2 + T_{1\bullet 2,3}^2$	Type A

TABLE I
TYPES DES DÉCOMPOSITIONS DU PROCÉDÉ À 3 VARIABLES

Li et al. [1] prouvent, en se basant sur les travaux d’Hawkins [28], que les décompositions de type A permettent un diagnostic moins précis que les décompositions de type B. De plus, ils prouvent également que dans le contexte du modèle linéaire gaussien, toutes les décompositions de type B convergent vers une unique décomposition que les auteurs nomment ”causation-based T^2 decomposition”. Nous la nommerons décomposition causale du T^2 . En effet, chaque décomposition de type B (dans le cas d’un modèle linéaire gaussien causal) converge vers la décomposition causale du T^2 décrite dans l’équation 6, où $PA(X_i)$ représentent les parents de la variable X_i sur le graphe causal.

$$T^2 = \sum_{i=1}^p T_{i\bullet PA(X_i)}^2 \quad (6)$$

Ainsi, la décomposition causale du T^2 de l’exemple de la figure 2 est la suivante : $T^2 = T_1^2 + T_{2\bullet 1}^2 + T_{3\bullet 1}^2$.

Suite à ces différentes démonstrations, les auteurs énoncent alors la procédure de détection et d’identification utilisant la nouvelle décomposition causale. Tout d’abord, un réseau bayésien linéaire gaussien est construit afin de représenter les relations causales entre les différentes variables du procédé. Suite à cela, le procédé est surveillé par une carte de contrôle du T^2 . Lors de la détection d’une situation hors-contrôle, le T^2 est décomposé par la décomposition causale de l’équation 6. Dans cette équation, chaque $T_{i\bullet PA(X_i)}^2$ est indépendant et, dans le cas où les paramètres du procédé sont connus, suit une distribution du χ^2 à un degré de liberté. On compare alors chaque $T_{i\bullet PA(X_i)}^2$ à la limite $\chi_{1,\alpha}^2$ représentant le quantile à la valeur α (taux de fausses alertes) de la distribution du χ^2 à un degré de liberté. Un $T_{i\bullet PA(X_i)}^2$ significatif (dépassant la limite de contrôle) implique alors que la variable X_i a probablement subi un saut de moyenne. La figure 3 représente le diagramme de surveillance du procédé par la méthode énoncée ci-dessus.

L’approche développée par Li et al. [1] exploite des seuils donnés par des quantiles de lois statistiques. Enfin, elle permet d’améliorer considérablement les performances par rapport à la méthode MYT, tout en demandant moins de calcul que celle-ci. Cependant, on remarque sur la figure 3 que cette méthode emploie divers outils : cartes de contrôle, réseaux bayésiens, calculs statistiques. Nous allons montrer, dans la section suivante, que la surveillance d’un procédé par la méthode de décomposition causale peut être entièrement effectuée par réseaux bayésiens. Ceci va

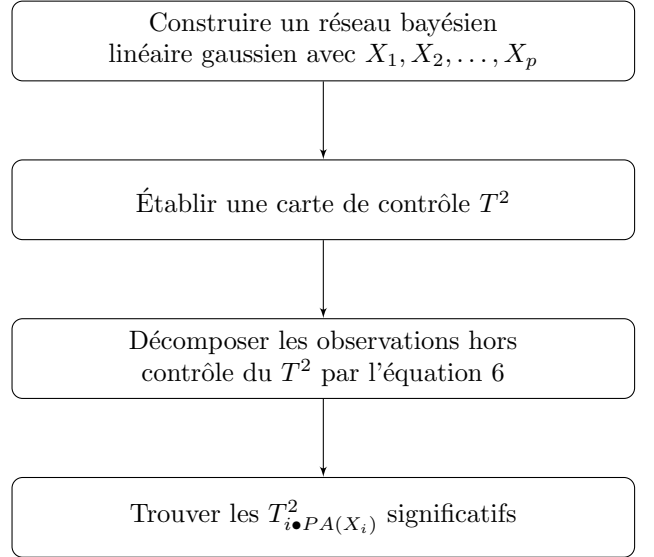


Fig. 3. Surveillance par la méthode de décomposition causale

notamment permettre la manipulation d’un seul et même outil pour la surveillance : un réseau bayésien.

IV. APPROCHE PROPOSÉE

A. Carte de contrôle par réseaux bayésiens

Lors de précédent travaux [2], [3], nous avons démontré qu’une carte de contrôle du T^2 [7] pouvait être modélisée par réseaux bayésiens. Pour cela, nous utilisons deux nœuds : un nœud multivariée gaussien \mathbf{X} représentant les données et un nœud bimodal E représentant l’état du procédé. Le nœud bimodal E possède les modalités suivantes : SC pour sous contrôle et HC pour hors-contrôle. En faisant l’hypothèse que $\boldsymbol{\mu}$ et $\boldsymbol{\Sigma}$ sont respectivement le vecteur cible et la matrice de variance-covariance du procédé, nous pouvons surveiller ce dernier par la règle suivante : si $p(SC|\mathbf{x}) < p(SC)$ alors le procédé est hors-contrôle. Ce réseau bayésien est représentée sur la figure 4, où nous précisons également les tables de probabilités conditionnelles associées à chaque nœud.

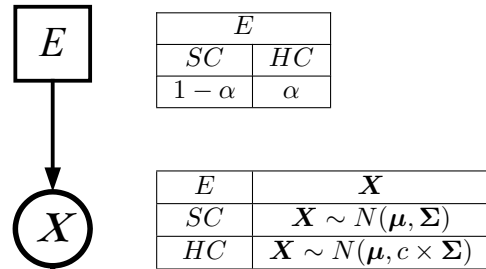


Fig. 4. Carte T^2 par réseaux bayésiens

Sur cette figure 4, nous pouvons observer qu’un coefficient c est impliqué dans la carte de contrôle par réseau bayésien. Ce coefficient est la racine (différente de 1) de l’équation suivante :

$$1 - c + \frac{pc}{LC} \ln(c) = 0 \quad (7)$$

où p est la dimension (nombre de variables) du système à surveiller et LC est la limite de contrôle de la carte T^2

équivalente. Dans de nombreux cas, LC est égale à $\chi_{\alpha,p}^2$, le quantile à la valeur α de la distribution du χ^2 à p degré de liberté [31]. Ainsi, α permet de régler le taux de fausses alertes de la carte de contrôle.

B. Amélioration de la proposition de Li et al.

La méthode proposée par Li et al. [1] permet, en se basant sur un réseau bayésien à nœuds gaussiens, de connaître les différents termes de la décomposition MYT à calculer. Pour le calcul des différents termes de la décomposition causale du T^2 , ainsi que pour les décisions associées (dépassement de limites), les auteurs n'utilisent pas leur réseau de façon optimale. En effet, les auteurs utilisent une carte de contrôle du T^2 extérieurement au réseau, alors que celle-ci peut se modéliser directement à l'intérieur du réseau (voir section précédente). De même, les auteurs calculent chaque $T_{i \bullet PA(X_i)}^2$ à l'extérieur du réseau, alors qu'il est possible d'effectuer ces calculs dans le réseau.

Nous proposons une extension à la méthode de Li et al. [1] permettant le calcul des différents $T_{i \bullet PA(X_i)}^2$ et des décisions associées à chacun d'entre eux. Le diagnostic par décomposition causale du T^2 , tout comme la décomposition MYT, est en fait une surveillance des variables régressées, au moyen de cartes de contrôle univariées. Dans la section précédente, nous avons démontré comment réaliser, par réseaux bayésiens, une carte de contrôle multivariée telle que la carte du T^2 de Hotelling. Or, une carte de contrôle univariée du type carte de contrôle de Shewhart n'est tout simplement qu'un cas particulier d'une carte de contrôle multivariée du type carte de T^2 de Hotelling. En effet, le calcul du T^2 est le suivant :

$$T^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad (8)$$

Or, dans le cas univarié, $\mathbf{x} = x$, $\boldsymbol{\mu} = \mu$ et $\boldsymbol{\Sigma} = \sigma^2$, ainsi l'équation 8 devient :

$$T^2 = \frac{(x - \mu)^2}{\sigma^2} \quad (9)$$

Dans ce cas univarié, la statistique T^2 suit une loi du χ^2 à un degré de liberté. Or, au vu des démonstrations de la section, ainsi que de leurs transpositions au domaine univarié, il est possible d'envisager une amélioration de la technique développée par Li et al. [1]. Nous proposons ici de suivre directement les différentes valeurs des $T_{i \bullet PA(X_i)}^2$ dans le réseau bayésien. Pour cela, nous rajoutons une variable discrète pour chaque nœud univarié du réseau bayésien. Si nous avons un graphe représentant un système à 3 variables (voir figure 2), nous obtenons alors un réseau avec six nœuds : 3 continus (univariés) et 3 discrets (bimodale), comme indiqué sur la figure 5.

Nous précisons ici que les variables continues n'ont pas obligatoirement besoin d'avoir été préalablement centrées et réduites. Les nœuds discrets rajoutés à la structure initiale du réseau (celle ne comprenant que les nœuds continus) nous permettent de réaliser directement l'identification des variables incriminées lors d'une situation hors-contrôle. Ces nœuds modélisent une carte de contrôle $T_{i \bullet PA(X_i)}^2$ permettant de conclure sur le statut de chaque variable. Nous rappelons que la modalité SC du nœud discret signifie sous contrôle, alors que la modalité HC signifie

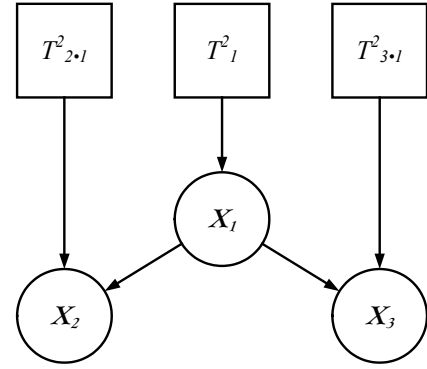


Fig. 5. Amélioration de la décomposition causale

hors-contrôle. La figure 6 détaille la table de probabilités conditionnelles associée à un nœud continu du réseau, ainsi que la table de probabilités a priori de son nœud discret associé.

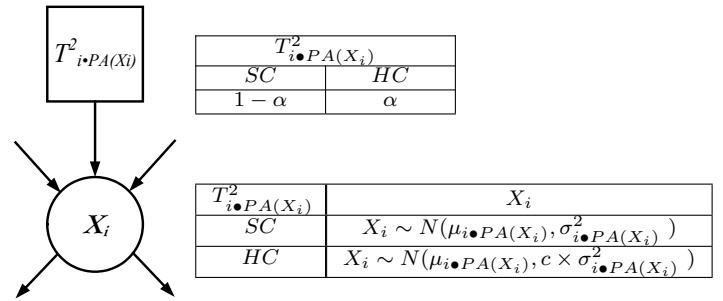


Fig. 6. Réseau bayésien similaire à la carte $T_{i \bullet PA(X_i)}^2$

Lorsqu'une faute est détectée dans le procédé, chaque nœud discret (représentant le statut d'une variable régressée) fournit une certaine probabilité que la variable soit sous contrôle. Les variables incriminées dans la faute du procédé sont les variables possédant une probabilité a posteriori inférieure à leur probabilité a priori. La personne chargée d'identifier physiquement la faute possède alors des indications très précieuses puisqu'elle connaît les variables du procédé sur lesquelles la faute a agi.

À la vue de ces propositions, nous pouvons dresser un réseau bayésien permettant, après une prise de décision sur les variables, de déterminer tout d'abord si le procédé est sous contrôle ou hors-contrôle (détection). Dans le cas du procédé hors contrôle, le réseau va également pouvoir nous fournir les variables responsables de cette situation (identification).

La figure 7 présente la forme générale de ce réseau bayésien pour un procédé à p variables. Ce réseau permet à lui seul d'effectuer toutes les phases de surveillance développées par Li et al. [1], c'est à dire qu'il permet l'application du diagramme de la figure 3.

V. APPLICATION

Afin de mieux appréhender la méthode proposée, nous avons simulé un procédé à trois variables tel que présenté sur les figures 2 et 5. Les paramètres ($\boldsymbol{\mu}$ et $\boldsymbol{\Sigma}$) de ce procédé lorsqu'il est sous contrôle sont les suivants :

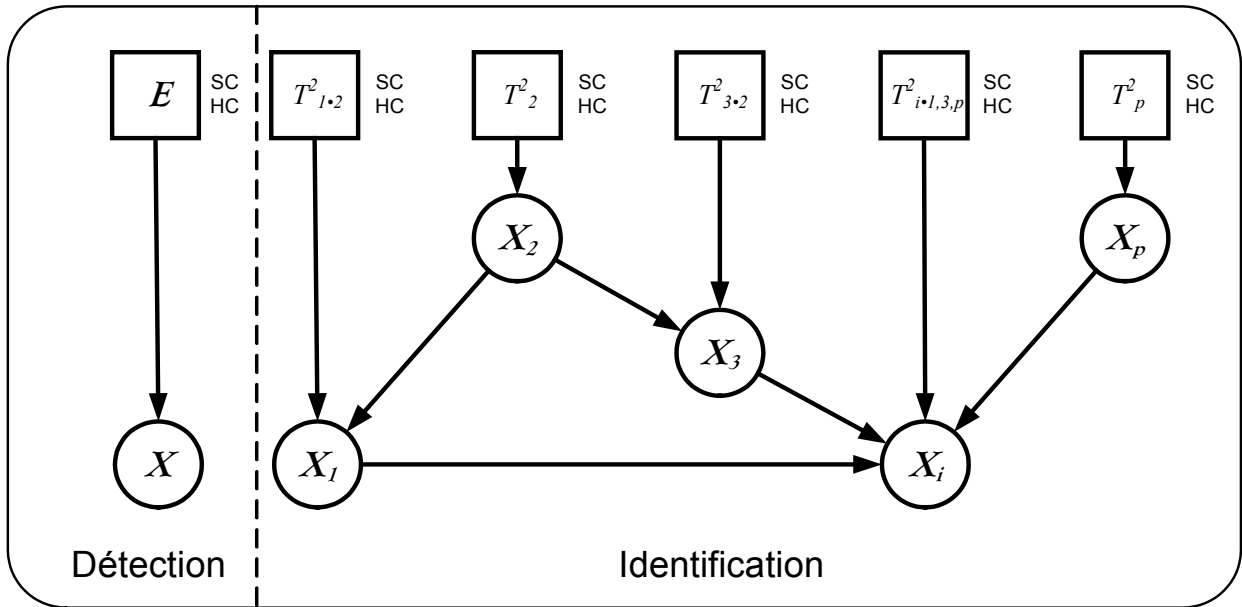


Fig. 7. Réseau bayésien pour la détection et l'identification

$$\mu = (0 \ 0 \ 0)$$

$$\Sigma = \begin{pmatrix} 1 & 0.8 & 0.6 \\ 0.8 & 1 & 0.1 \\ 0.6 & 0.1 & 1 \end{pmatrix}$$

Nous avons simulé 100 observations sous contrôle, puis 50 observations hors contrôle. Ces observations hors contrôle sont un saut de moyenne d'amplitude 2 sur la variable X_2 . L'objectif du réseau est donc de détecter qu'une faute est apparue à partir de l'observation 101, et que la variable X_2 est impliquée dans cette faute.

La première action effectuée est la construction du réseau conditionnel gaussien. Pour cela, l'algorithme PC [32] est utilisé pour la construction de la structure du réseau. Pour l'algorithme PC, nous avons employé un test d'indépendance conditionnelle par transformé en z [33]. On ajoute alors tous les nœuds de contrôle de paramètres, c'est à dire tous les nœuds $T_{i \bullet PA(X_i)}^2$, avec un taux de fausses alertes de 5%. On ajoute également la modélisation de la carte de contrôle du T^2 , permettant la détection par le réseau, avec un taux de fausses alarmes de 1%. Suite à cela, les paramètres de chaque nœud sont calculés.

Suite à l'apprentissage de la structure et des paramètres du réseau, nous présentons les 150 observations pour lesquelles nous voulons obtenir un résultat. On observe alors les probabilités a posteriori du nœud E . Pour une observation donnée, si cette probabilité est inférieure à 99% alors le procédé est hors contrôle. Dans ce cas, on observe les valeurs des probabilités a posteriori des nœuds $T_{i \bullet PA(X_i)}^2$. Les variables possédant une probabilité inférieure à 95% sont alors impliquées dans la situation de hors contrôle. Les résultats (les probabilités à posteriori) du réseau sont donnés sur la figure 8.

Sur cette figure, on peut remarquer que le réseau ne détecte pas de faute sur les 100 premières observations. Cependant, le réseau détecte une situation hors contrôle

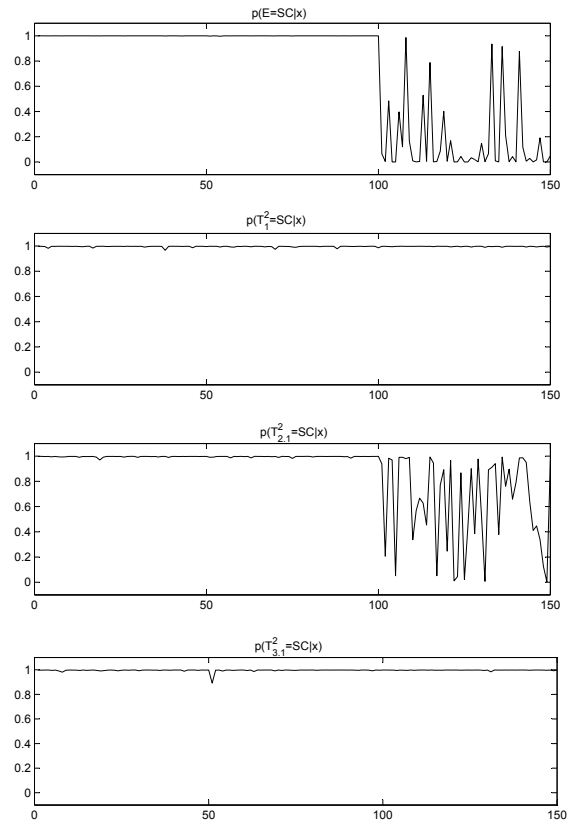


Fig. 8. Probabilités à posteriori des nœuds

dès l'observation 101. De plus, le réseau nous signale que la variable X_2 est impliquée dans cette faute.

VI. CONCLUSIONS ET PERSPECTIVES

Dans cet article nous avons présenté une approche permettant d'effectuer la détection de fautes dans un procédé multivarié, ainsi que l'identification des différentes variables impliquées lors de ces fautes. Cette approche se

base sur un réseau bayésien. Il s'agit de la combinaison de précédents travaux [2], [3] permettant la modélisation de cartes de contrôle multivariées par réseau bayésien, ainsi que de l'amélioration de la technique de décomposition causale du T^2 proposée par Li et al. [1]. Ainsi, la perspective intéressante qu'entraîne cette méthode est la mise au point d'un réseau bayésien complet permettant à la fois la détection, l'identification et le diagnostic de faute dans un procédé multivarié.

RÉFÉRENCES

- [1] Jing Li, Jionghua Jin, et Jianjun Shi. Causation-based t^2 decomposition for multivariate process monitoring and diagnosis. *Journal of Quality Technology*, 40(1) :46–58, 2008.
- [2] Sylvain Verron, Teodor Tiplica, et Abdessamad Kobi. Multivariate control charts with a bayesian network. *4th International Conference on Informatics in Control, Automation and Robotics (ICINCO)*, 2007.
- [3] Sylvain Verron. *Diagnostic et surveillance des processus complexes par réseaux bayésiens*. PhD thesis, University of Angers, 2007.
- [4] Leo H. Chiang, Evan L. Russell, et Richard D. Braatz. *Fault detection and diagnosis in industrial systems*. New York : Springer-Verlag, 2001.
- [5] Walter A. Shewhart. *Economic control of quality of manufactured product*. New York : D. Van Nostrand Co., 1931.
- [6] S. W. Roberts. Control chart tests based on geometric moving averages. *Technometrics*, 1(3) :239–250, Août 1959.
- [7] Harold Hotelling. Multivariate quality control. *Techniques of Statistical Analysis*, :111–184, 1947.
- [8] J.A. Westerhuis, S.P. Gurden, et A.K. Smilde. Standardized q-statistic for improved sensitivity in the monitoring of residuals in mspc. *Journal of Chemometrics*, 14(4) :335–349, 2000.
- [9] Edward J. Jackson. Multivariate quality control. *Communication Statistics - Theory and Methods*, 14 :2657 – 2688, 1985.
- [10] Bhavik R. Bakshi. Multiscale PCA with application to multivariate statistical process monitoring. *AIChE Journal*, 44 :1596–1610, 1998.
- [11] M. Kano, K. Nagao, S. Hasebe, I. Hashimoto, H. Ohno, R. Strauss, et B.R. Bakshi. Comparison of multivariate statistical process monitoring methods with applications to the eastman challenge problem. *Computers and Chemical Engineering*, 26(2) :161–174, 2002.
- [12] J.F. MacGregor et T. Kourti. Statistical process control of multivariate processes. *Control Engineering Practice*, 3(3) :403–414, 1995.
- [13] B.M. Wise et N.B. Gallagher. The process chemometrics approach to process monitoring and fault detection. *Journal of Process Control*, 6(6) :329–348, 1996.
- [14] Teodor Tiplica, Abdessamad Kobi, et Alain Barreau. Synthèse et comparaison des méthodes pour la maîtrise statistique des processus multivariés. *Actes du congrès QUALITA*, pages 134–142, Annecy, France, 2001.
- [15] Robert L. Mason, Nola D. Tracy, et John C. Young. Decomposition of T^2 for multivariate control chart interpretation. *Journal of Quality Technology*, 27(2) :99–108, 1995.
- [16] R.L. Mason, N.D. Tracy, et J.C. Young. A practical approach for interpreting multivariate t^2 control chart signals. *Journal of Quality Technology*, 29(4) :396–406, 1997.
- [17] Finn V. Jensen. *An introduction to Bayesian Networks*. Taylor and Francis, London, United Kingdom, 1996.
- [18] Patrick Naim, Pierre-Henri Wuillemin, Philippe Leray, Olivier Pourret, et Anna Becker. *Réseaux bayésiens - 2ème édition*. Eyrolles, 2004.
- [19] Pat Langley et Stephanie Sage. Induction of selective bayesian classifiers. *In Proceedings of the Tenth Conference on Uncertainty in Artificial Intelligence*, 1994.
- [20] N. Friedman, D. Geiger, et M. Goldszmidt. Bayesian network classifiers. *Machine Learning*, 29(2-3) :131–163, 1997.
- [21] F. Pernkopf. Bayesian network classifiers versus selective k-nn classifier. *Pattern Recognition*, 38(1) :1–10, 2005.
- [22] Pedro Domingos et Michael J. Pazzani. Beyond independence : Conditions for the optimality of the simple bayesian classifier. *International Conference on Machine Learning*, 1996.
- [23] C. Chow et C. Liu. Approximating discrete probability distributions with dependence trees. *Information Theory, IEEE Transactions on*, 14(3) :462–467, 1968.
- [24] Igor Kononenko. Semi-naive bayesian classifier. *EWSL-91 : Proceedings of the European working session on learning on Machine learning*, pages 206–219, 1991.
- [25] M. Pazzani. Searching for dependencies in bayesian classifiers. *Learning from Data Artificial Intelligence and Statistics*, 5 :239–248, 1997.
- [26] Meng Koon Chua et Douglas C. Montgomery. Investigation and characterization of a control scheme for multivariate quality control. *Quality and Reliability Engineering International*, 8 :37–44, 1992.
- [27] Necip Doganaksoy, Frederick .W. Faltin, et William T. Tucker. Identification of out of control quality characteristics in a multivariate manufacturing environment. *Communications in Statistics - Theory and Methods*, 20(9) :2775–2790, 1991.
- [28] Douglas M. Hawkins. Regression adjustment for variables in multivariate quality control. *Journal of Quality Technology*, 25(3) :170–182, 1993.
- [29] Teodor Tiplica, Abdessamad Kobi, et Alain Barreau. Optimisation et maîtrise des processus multivariés. la méthode fnad. *Journal Européen des Systèmes Automatisés*, 37(4) :477–500, 2003.
- [30] Teodor Tiplica. *Contribution à la Maîtrise Statistique des Processus Industriels Multivariés*. PhD thesis, ISTIA, 2002.
- [31] Douglas C. Montgomery. *Introduction to Statistical Quality Control, Third Edition*. John Wiley and Sons, 1997.
- [32] Peter Spirtes, Clark Glymour, et Richard Scheines. *Causation, prediction, and search*. Springer-Verlag, 1993.
- [33] M. Kalisch et P. Buhlmann. Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *Journal of Machine Learning Research*, 8 :613–636, 2007.