

## MONITORIZACIÓN DE PROCESOS POR LOTES MEDIANTE PCA MULTIFASE

José Camacho <sup>\*,1</sup> Jesús Picó <sup>\*</sup>

*\* Departamento de Ingeniería de Sistemas y Automática  
Universidad Politécnica de Valencia*

Resumen: Las técnicas estadísticas, como el Análisis por Componentes Principales (PCA), son ampliamente utilizadas para la monitorización de procesos por lotes. Sin embargo, estas técnicas se basan en la existencia de un modelo lineal que aproxime bien los datos, por lo que su rendimiento se reduce al modelar conjuntos de datos de naturaleza no lineal. Una solución es la división del modelo en varios submodelos, uno por cada segmento del lote que pueda ser bien aproximado por un modelo lineal. En este artículo se propone un esquema de monitorización basado en la detección automática de los segmentos lineales a lo largo del lote. Cada segmento lineal será modelado independientemente. Finalmente se presenta una modificación de las tablas de monitorización para incluir conjuntamente los múltiples modelos.  
*Copyright ©2006 CEA-IFAC*

Palabras Clave: Control Estadístico Multivariable de Procesos, PCA Multifase, Monitorización, Procesos por Lotes, Modelos locales.

### 1. INTRODUCCIÓN

Los procesos por lotes se caracterizan por su flexibilidad para la producción, en una misma línea, de distintos productos de poco volumen pero elevado valor. Estos procesos han sido utilizados en la industria farmacéutica, bioquímica, alimentaria, cerámica, metalúrgica y electrónica, entre otros.

Con carácter general, los procesos por lotes se pueden definir como procesos de duración finita, consistentes en la repetición continuada de tres pasos (Nomikos and MacGregor, 1995): carga de la bandeja<sup>2</sup> con una receta específica, procesado bajo condiciones controladas y descarga de la bandeja. Una vez que el lote ha sido procesado, se pueden obtener un conjunto de medidas

que caracterizan la calidad de éste. Dependiendo de la disponibilidad y coste de sensores para la obtención de estas medidas, su evolución puede observarse a lo largo del procesamiento, pueden estar disponibles únicamente al final de cada lote o bien han de ser obtenidas con posterioridad en laboratorio. En este último caso, puede que no sea posible disponer de las medidas de calidad de un lote al inicio del procesamiento del siguiente.

En la literatura se ha prestado menos atención a este tipo de procesos que a los procesos continuos. Las técnicas utilizadas para monitorizar y controlar los procesos continuos no son directamente aplicables a los procesos por lotes (Trelea *et al.*, 1997), ya que estos últimos presentan una dimensión adicional con la que tratar: el número de lote (Figura 1). Mientras que en los procesos continuos se busca mantener un punto de operación o estado estable, el buen funcionamiento en un proceso por lotes se obtiene cuando el lote bajo procesamiento describe una trayectoria cercana a

<sup>1</sup> Parcialmente financiado a través del programa de becas de Formación de Profesorado Universitario (FPU), Secretaría de Estado de Educación y Universidades

<sup>2</sup> Recipiente en sentido amplio, en el que se procesa la unidad del producto

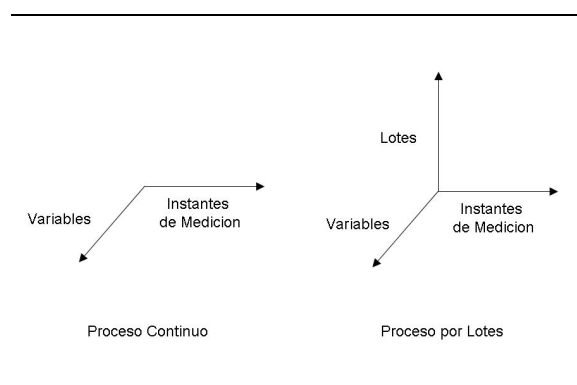


Figura 1. Naturaleza bidimensional de los procesos continuos frente a la tridimensional de los procesos de lotes

una trayectoria nominal dada. En la Figura 2 se representan las trayectorias sin alinear de 50 lotes para 3 de las variables de un proceso por lotes.

El principal objetivo del control de procesos por lotes es la fabricación de productos acordes a ciertas especificaciones sobre las medidas de calidad, maximizando la utilización de la maquinaria disponible (Edgar, 2004). El sistema de monitorización ha de establecer un modelo de las Condiciones de Operación Normal (NOC) a partir del cual realizar tareas como la detección de fallos. De esta manera, cuando un lote presenta un comportamiento distinto al conjunto de lotes que cumplen las especificaciones, es catalogado como erróneo y, si es posible actuar durante la ejecución del lote, se tomarán las acciones correctoras pertinentes.

La automatización de la industria ha llevado a la aparición de amplias bases de datos que permiten la generación de modelos de caja negra. En los procesos por lotes, los datos de la base histórica son de carácter tridimensional (Figura 1). Este hecho, conjuntamente con que las medidas de calidad puedan no estar disponibles en línea, hace el control especialmente complejo.

Las bases históricas a menudo están formadas por datos altamente correlacionados y de baja relación señal ruido (Kourti, 2002). Esto hace que la información no aparezca explícitamente en ellos y tengan que utilizarse técnicas que permitan su extracción (Singhal and Seborg, 2002). Además, estos datos presentan una serie de características específicas que hacen necesario el preprocesamiento para que algunas de las técnicas de extracción de información puedan ser utilizadas en su análisis (Ündey and Çinar, 2002):

- Dinámica no lineal y variable en el tiempo.
- Longitud distinta de los lotes.
- Presencia de ruido, datos colineales, valores perdidos y muestras alejadas de la distribución de datos.
- Diferentes magnitudes y varianzas entre las variables.

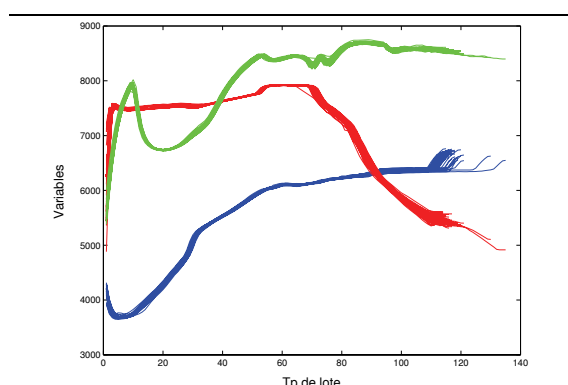


Figura 2. Trayectoria de 3 de las variables de un proceso por lotes.

En este contexto, las técnicas basadas en la búsqueda de variables latentes son las más frecuentemente aplicadas para la generación de modelos empíricos. Estas técnicas permiten mejorar el conocimiento sobre el proceso (Kosanovich *et al.*, 1996). Sin embargo, en su forma tradicional, este tipo de técnicas únicamente son capaces de capturar relaciones lineales entre variables. Una estrategia para afrontar la no linealidad es la identificación de modelos locales.

En este artículo se presenta un algoritmo para la identificación de fases en el procesamiento de un lote, que puedan ser convenientemente modeladas a través de un modelo lineal. Adicionalmente, las tablas de monitorización tradicionales son extendidas para el uso de modelos multifase. En (Camacho and Picó, 2006) se presenta una versión extendida de este algoritmo y una comparativa con otras técnicas multifase.

El artículo se organiza de la siguiente manera: en la sección 2 se introduce el modelado de procesos por lotes a partir de técnicas estadísticas, en concreto del Análisis por Componentes Principales (PCA); en la sección 3 se describe el algoritmo propuesto para la detección de segmentos lineales en los datos; en la sección 4, se presentan los dos conjuntos de datos utilizados en las pruebas; la sección 5 ofrece resultados de la aplicación del algoritmo sobre los conjuntos de datos; en la sección 6 se presentan las modificaciones sobre las tablas de monitorización tradicionales, para el uso de modelos multifase; finalmente, en la sección 7 se extraen conclusiones y se avanzan líneas futuras de trabajo.

## 2. TÉCNICAS ESTADÍSTICAS MULTIVARIABLES PARA MODELAR PROCESOS DE LOTES

El carácter poco explícito de la información contenida en las bases históricas hace que la aplicación de técnicas estadísticas sea especialmente adecuada. El modelado estadístico consiste en ex-

traer de la base histórica un conjunto de lotes de resultado satisfactorio para representar las NOC, y analizar este conjunto para modelarlo. En este punto es relevante resaltar que el concepto de modelo estadístico no es análogo al de modelo dinámico, al cual está habituado el ingeniero de control.

Las técnicas más destacadas en la literatura son el Análisis de Componentes Principales (PCA) y la Proyección sobre Estructuras Latentes (PLS), también llamada Mínimos Cuadrados Parciales. PCA y PLS son técnicas lineales que requieren el preprocesamiento de los datos. En general, se asume que la eliminación de la trayectoria media del NOC elimina a su vez la mayor parte de la no linealidad en el lote (Nomikos and MacGregor, 1995), y por tanto son directamente aplicables técnicas lineales de modelado (Kosanovich *et al.*, 1996).

Para la creación de un modelo estadístico que permita establecer cuándo una variación con respecto a la trayectoria nominal ha de ser corregida, es necesario que todos los lotes sean de igual longitud y deseable que presenten las mismas características en los mismos puntos de medición. Una manera sencilla de conseguir esta equalización es el uso de una variable indicadora (Nomikos and MacGregor, 1994). Esta variable actúa como el reloj interno del proceso de fabricación de un lote. Si las mediciones se han realizado en escala temporal, las correspondientes mediciones en la escala de la variable indicadora se pueden obtener por interpolación. Esta solución, sin embargo, no es aplicable en todos los casos, ya que es imprescindible la existencia de una variable de estas características.

Otra posibilidad es equalizar los datos haciendo corresponder características comunes en la trayectoria de los lotes. Éste es el caso del Ajuste de Tiempo Dinámico (DTW), que permite equalizar cada lote según una trayectoria nominal dada. La aplicación de esta técnica a la monitorización de los procesos de lotes ofrecida en (Kassidas *et al.*, 1998) puede empeorar la detección de fallos, como indican los propios autores.

En cuanto al carácter ruidoso de los datos y las diferencias entre unidades o varianzas en las variables, es recomendable filtrar y escalar previo al análisis estadístico.

### 2.1 Análisis de Componentes Principales

PCA se sustenta en la idea de que, en una base de datos formada por gran cantidad de variables muy correlacionadas entre sí, la dimensión real del espacio donde se ubican los valores de estas variables es pequeña (Kourti, 2002). Esto permite

sustituir el conjunto de variables original por un número mucho menor de combinaciones de éstas, las variables latentes, con poca pérdida de información. Esta reducción hace posible que se pueda utilizar de forma efectiva la información contenida en las bases históricas en tareas como monitorizar o controlar un proceso.

PCA está definido para el análisis de matrices bidimensionales de la forma  $\mathbf{X}(I \times J)$ , donde se almacena un conjunto de  $I$  observaciones sobre  $J$  variables. Permite obtener, dentro del espacio de todas las variables, el subespacio que encierre la mayor parte de la variabilidad de las muestras. Esta técnica establece los ejes del espacio en las direcciones de mayor variabilidad, para descartar las direcciones que explican menos varianza y por tanto contienen poca información. El algoritmo iterativo NIPALS (Wold and Lyttkens, 1969) (ver Apéndice) permite calcular estas direcciones ordenadamente, de manera que los ejes a descartar no han de ser calculados.

El modelo obtenido con PCA se corresponde con la ecuación (1).

$$\mathbf{X} = \mathbf{TP}' + \mathbf{E} \quad (1)$$

donde la matriz de puntuaciones (*scores*)  $\mathbf{T}(I \times R)$  es el conjunto de las proyecciones de las  $I$  observaciones en el espacio  $R$ -dimensional, con  $R$  menor o igual al número de variables  $J$ , la matriz de cargas (*loadings*)  $\mathbf{P}(J \times R)$  contiene en sus columnas los  $R$  autovectores de la matriz  $\mathbf{X}'\mathbf{X}$  con mayor autovalor asociado, y la matriz de residuos  $\mathbf{E}(I \times J)$  contiene la información descartada por el modelo.

Cada una de las  $R$  variables latentes, o componentes principales, resulta de la combinación de las  $J$  variables iniciales según se especifica en un autovector de  $\mathbf{P}$ . La ecuación (2) realiza el cálculo del valor (o puntuación) de una nueva observación en el espacio reducido de las  $R$  variables latentes.

$$\mathbf{t} = \mathbf{x}_{\text{new}}\mathbf{P} \quad (2)$$

La técnica más utilizada para calcular el número de componentes principales que permite representar adecuadamente el comportamiento del NOC es la validación cruzada. El procedimiento habitual es calcular el error cuadrático de predicción o PRESS para un cierto número de componentes principales (por ejemplo, de 1 a 10). A partir de la forma del PRESS o usando algún índice estadístico (Wold, 1978) se determina el número idóneo de componentes.

### 2.2 Aplicación de PCA a los Procesos por Lotes

PCA permite el análisis de datos bidimensionales. Como se ha comentado con anterioridad, los datos

obtenidos durante el funcionamiento de una planta de procesamiento por lotes son de naturaleza tridimensional. Una solución es la transformación de la matriz  $\underline{\mathbf{X}}(I \times J \times K)$ , con  $I$  lotes,  $J$  variables y  $K$  instantes de medición, en una matriz bidimensional (Nomikos and MacGregor, 1994)(Wold *et al.*, 1987). De esta propuesta nace la técnica PCA multidireccional o MPCA (Multiway-PCA). MPCA aplica tres pasos para realizar el análisis.

En primer lugar se despliega la matriz en la dirección de una dimensión, intercalando las otras dos dimensiones. Existen 6 formas de realizar este despliegue, ya que hay tres direcciones de despliegue y la posibilidad de trasponer o no la matriz resultante. El análisis de las 6 matrices resultantes explica un tipo distinto de variabilidad (Kosanovich *et al.*, 1996).

Nomikos y MacGregor propusieron en (Nomikos and MacGregor, 1994) el despliegue de la matriz en la dimensión de los lotes (3), mostrada en Figura 3(a), que permite analizar la variabilidad entre lotes y es la forma ideal para llevar a cabo tareas de monitorización, debido al preprocesamiento realizado (Westerhuis *et al.*, 1999). Tras este despliegue, en la matriz bidimensional resultante cada lote aparece como una única observación y el número de variables es aumentado  $K$  veces. La matriz resultante de varianzas-covarianzas (Figura 3(b)), sobre la cual se aplica PCA, permite estudiar las relaciones entre todas las variables en el mismo (submatrices  $\{VC_1 \dots VC_K\}$ ) y distintos instantes temporales (resto de la matriz).

$$\underline{\mathbf{X}}(I \times J \times K) \Rightarrow \mathbf{X}(I \times JK) \quad (3)$$

El segundo paso es el análisis de la matriz desplegada utilizando PCA. Una vez que se ha obtenido la matriz de cargas  $\mathbf{P}(JK \times R)$ , se vuelve a plegar en  $\underline{\mathbf{P}}(J \times K \times R)$  en un tercer paso. El modelo final aparece en la ecuación (4).

$$\hat{x}_{ijk} = \sum_{r=1}^R t_{ir} p_{jkr} \quad (4)$$

donde  $\hat{x}_{ijk}$  es la predicción del valor de la variable  $j$  en el punto de medición  $k$  del lote  $i$ ,  $t_{ir} \in \mathbf{T}(I \times R)$  es la puntuación del lote  $i$  en el componente principal  $r$ ,  $p_{jkr} \in \underline{\mathbf{P}}(J \times K \times R)$  es la carga de la variable  $j$  en el punto de medición  $k$  en el componente principal  $r$  y  $\mathbf{X} = \hat{\mathbf{X}} + \underline{\mathbf{E}}$ , con  $\mathbf{E}(I \times J \times K)$  la matriz de residuos.

### 3. GENERACIÓN AUTOMÁTICA DE MODELOS PCA MULTIFASE

En el modelado de procesos de lotes, uno de los mayores inconvenientes para la aplicación de

MPCA es su carácter lineal. Aunque el procedimiento común para el análisis de datos tridimensionales es extraer la trayectoria media en el tiempo, los datos resultantes siguen presentando no-linealidades que complican la aplicación de esta técnica (Rotem *et al.*, 2000). Para afrontar las no linealidades, se puede utilizar alguna de estas dos estrategias:

- Modificar PCA o los datos para que se puedan extraer variables latentes que representen combinaciones no lineales de las variables originales. Un ejemplo es el PCA no lineal propuesto en (Dong and McAvoy, 1996).
- Dividir el lote en segmentos, donde la estructura de correlación pueda ser aproximada con mayor precisión por variables latentes que representen combinaciones lineales (Técnicas multifase). De cada segmento se obtiene un modelo local.

El empleo de modelos locales ante la no linealidad del proceso por lotes no responde a la misma idea subyacente que en los procesos continuos. Un lote, a lo largo de su procesamiento, pasa por distintas unidades de proceso, distintas reacciones químicas, etc., cada una de ellas con una dinámica distinta. El proceso puede ser modelado en conjunto, o bien utilizando un modelo por cada una de estas fases. Estos modelos son locales a un intervalo del avance del lote. En los procesos continuos, los modelos locales aproximan una función no lineal alrededor de un punto de operación (Díez *et al.*, 2004). Son, pues, locales a este punto.

Desde el punto de vista estadístico, los modelos locales representan un subespacio del espacio total de variables, donde las muestras puedan ser bien aproximadas por un modelo lineal. Al desplegar la matriz tridimensional de datos en una matriz de dos dimensiones (3), cada variable original formará tantas variables como puntos de medición haya, de manera que el número de variables final es muy elevado (con 3 variables y 100 puntos de medición, el número de variables final será 300). Si dividimos el lote en intervalos, las variables de cada intervalo conformarán un espacio menor, facilitando la aproximación lineal.

El principal problema de las técnicas multifase es que requieren de conocimiento experto del proceso. Es decir, hay que saber cuándo el proceso es no lineal y hay necesidad de aplicarlas. Además, es necesario establecer la ubicación de la división en múltiples modelos. En (Louwerse and Smilde, 2000) se propone dividir el modelo en intervalos regulares y calcular los submodelos con las muestras desde el principio del lote hasta cada una de las divisiones. Esta subdivisión no responde a la naturaleza específica de cada proceso. En (Ündey and Çinar, 2002) se propone crear un modelo por cada unidad física de pro-

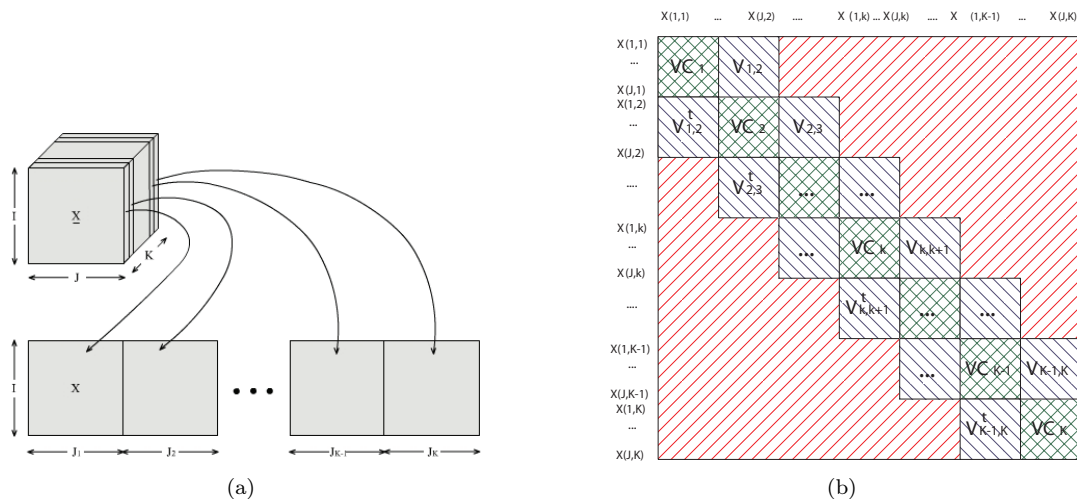


Figura 3. Despliegue propuesto por Nomikos y MacGregor (a) y matriz de varianzas-covarianzas resultante (b).

ceso (lo que llaman *stage*), y dentro de éstas, uno por cada subfase (lo que llaman *phase*, por ejemplo, distintas reacciones químicas dentro de cada unidad). Sin embargo, es necesario conocer la existencia y ubicación de estas subfases. Siguiendo una nomenclatura parecida, a partir de ahora llamaremos fases a los distintos pasos diferenciados por los que pasa un proceso por lotes y subfases a los segmentos del lote bien representados por un modelo lineal tras el preprocesamiento, que es lo que queremos encontrar.

Este artículo propone un algoritmo para la subdivisión del modelo completo del proceso en submodelos donde la estructura de correlación sea bien representada por un modelo lineal. Este algoritmo se definirá recursivo. De esta manera, cada vez que se proponga una separación, su conveniencia será comprobada y, si el resultado es positivo, se ejecutará recursivamente sobre las partes resultantes.

La estructura *Top-Down* del algoritmo y sus características específicas (argumentos de entrada y medida de bondad de una subdivisión) son completamente novedosas en la aplicación que nos ocupa. Con una visión radicalmente opuesta, las recientes propuestas en (Abonyi *et al.*, 2005) y (Lu *et al.*, 2004) se basan en técnicas de clustering. Estas técnicas no tratan la posición temporal de las muestras de forma implícita, como en el algoritmo propuesto. Esto lleva a que esta información se introduzca de forma explícita o a posteriori, modificando el resultado, lo que puede producir dificultades en la segmentación para algunos procesos.

### 3.1 Medida de comparación entre modelos para la subdivisión

Una medida apropiada para la comparación de modelos de iguales dimensiones es la varianza explicada. Esta medida ya ha sido utilizada para justificar la bondad del modelo multifase en el tipo de problemas al que nos enfrentamos (Ündey and Çinar, 2002)(Lu *et al.*, 2004).

El estudio de la varianza explicada nos permitirá evaluar la bondad de una determinada división. Para efectuar la comparación, se obtendrá el mismo número de componentes principales tanto para el modelo previo a la división como para el modelo dividido. De esta forma, el tamaño de ambos modelos será el mismo. Esta filosofía hace que todos los submodelos que se generen tengan el mismo número de componentes principales, ya que la división del lote en subfases está justificada para un valor concreto de componentes.

Los siguientes argumentos justifican las características de la medida:

- Al dividir en submodelos la varianza explicada aumenta, ya que cada modelo representa menor cantidad de datos y se ajusta mejor a éstos. Una mejora pequeña en la varianza explicada puede indicar que la separación en submodelos no está justificada. Por este motivo se establece un umbral como parámetro de entrada en el algoritmo.
- El umbral se fijará sobre la reducción relativa de varianza no explicada (residuos). Esta definición del umbral hace que el número de divisiones se mantenga constante en alto grado al variar el número de componentes principales, lo que no ocurre con otras definiciones estudiadas.

- c) La mejora de la varianza explicada es calculada sobre el mismo conjunto de muestras utilizado para la generación del modelo. Existen alternativas como la validación cruzada o el uso de un conjunto de test. La validación cruzada presenta mayor carga computacional (del orden de 60 veces más), ofreciendo una mejora muy reducida. De esta forma, únicamente se utilizará la validación cruzada para validar la mejor división encontrada, pero no en la búsqueda de ésta. Si se utiliza un conjunto de test, la división puede llegar a ser muy dependiente de éste. La elección del conjunto de test, por tanto, ha de realizarse con conocimiento experto del proceso. Precisamente el objetivo de este artículo es realizar el análisis sin disponer de este conocimiento.

### 3.2 Algoritmo PCA Multifase

El algoritmo PCA Multifase (Multi-Phase PCA, MPPCA), en cada recursión, realiza una búsqueda exhaustiva sobre todas las posibilidades a partir de un conjunto de parámetros de entrada, a saber:

- La matriz tridimensional con los datos (completa o una subdivisión en recursiones posteriores).
- El número de variables latentes de los modelos.
- El tamaño mínimo de submodelo, como porcentaje del tamaño total del lote.
- El umbral establecido o mejora mínima para que una división sea aceptada.

El algoritmo realiza los siguientes pasos:

- i. Despliegue de la matriz tridimensional.
- ii. Preprocesado de datos.
- iii. PCA del modelo completo.
- iv. Para cada instante de medición ( $K$ )
  - iv.1 Si la división genera dos modelos mayores al tamaño mínimo definido.
    - iv.11 PCA de los dos submodelos.
    - iv.12 Cálculo de la mejora.
    - iv.13 Si el resultado es el mejor hasta el momento. Actualiza los datos de la mejor subdivisión.
- v. Cálculo de la varianza explicada por validación cruzada de los modelos iii y iv. Acabar si no se supera el umbral de mejora definido.
- vi. Aceptar submodelos y ejecutar pasos iii a vi para cada submodelo.

El desplegado y preprocesamiento utilizado es el propuesto por Nomikos y MacGregor (Nomikos and MacGregor, 1994).

El algoritmo parte del modelo completo y busca la mejor división posible (en términos de varianza explicada). En segundo lugar, esta división es

comparada (usando validación cruzada) con el modelo sin dividir. Si la división provoca una mejora superior al umbral, ésta es aceptada y se prosigue de forma recursiva por los submodelos obtenidos.

Para determinar el número de componentes principales del modelo MPPCA se sigue el mismo procedimiento explicado en la sección 2.1.

## 4. CONJUNTOS DE DATOS DE MUESTRA

El conjunto de pruebas presentadas en este artículo ha sido realizado sobre dos conjuntos de muestras. El primero pertenece al proceso de polimerización de nylon-6,6 utilizado en (Kosanovich *et al.*, 1996). En el artículo anterior se propone la división en dos submodelos. El segundo es el NOC de un proceso de serigrafía de tarjetas electrónicas que aparece en (Wise *et al.*, 1999). Este conjunto de datos presenta distintos puntos de operación de la planta. Como nuestro objetivo no es aplicar ninguna técnica adaptativa, se han seleccionado los 30 primeros lotes, que pertenecen al mismo punto de operación.

El primer conjunto de datos no aparece alineado, pero cada medida temporal incorpora información sobre su pertenencia a una de las fases físicas del proceso. Los lotes se alinean a partir de esta información por interpolación lineal. El conjunto final es una matriz tridimensional de 50 lotes, 9 variables y 116 puntos de medición. El segundo conjunto de datos aparece previamente alineado y consta de 30 lotes, 12 variables y 80 puntos de medición.

## 5. APLICACIÓN DEL ALGORITMO

Para obtener un conjunto de muestras que nos permitan determinar cómo se comporta el algoritmo MPPCA al variar los parámetros de entrada, se han utilizado los siguientes valores:

- El número de componentes principales variará entre 1 y 10.
- El tamaño mínimo de submodelo será 1/3, 1/4, 1/5, 1/6, 1/7 y 1/8 del tamaño de lote. Estos valores limitan el número de submodelos posibles a 3, 4, 5, 6, 7 y 8, respectivamente.
- El rango de umbrales utilizado irá del 1% al 20% de mejora, en pasos de 1.

Como medida de la bondad de un modelo se ha utilizado la varianza explicada calculada a través de validación cruzada, con el procedimiento "deja uno fuera". La varianza explicada permite evaluar la adecuación del modelo a los datos y su cálculo

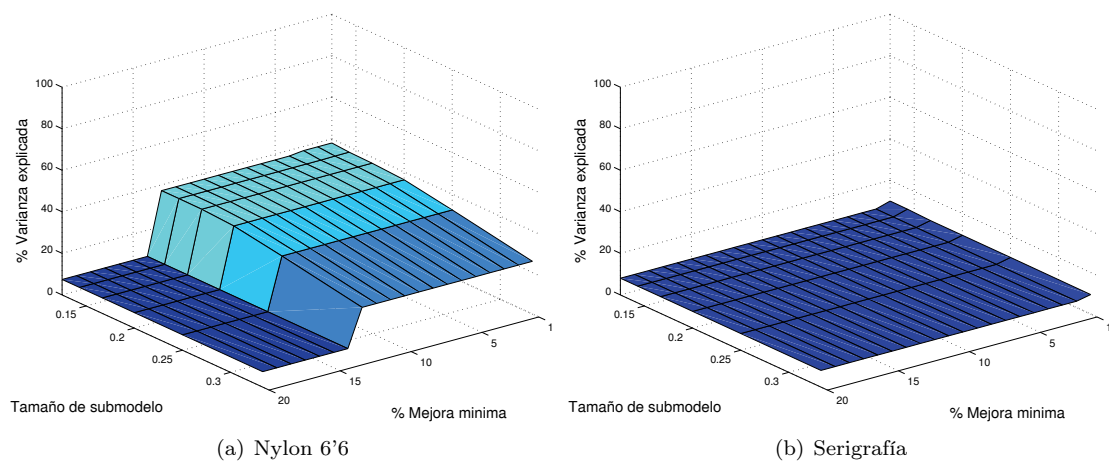


Figura 4. Varianza explicada tras la separación para 1 componente principal.

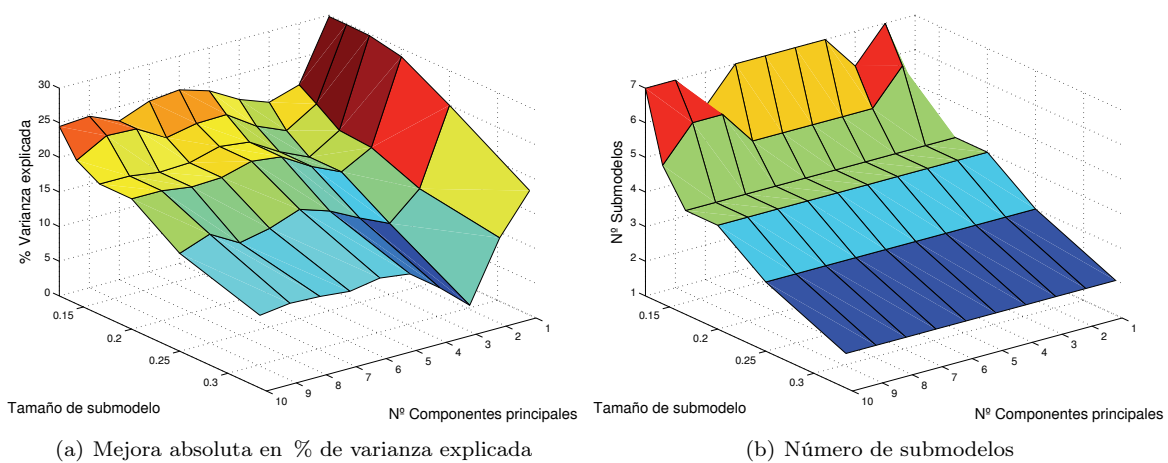


Figura 5. Resultados de MPPCA con umbral del 7% sobre el proceso nylon 6'6.

por validación cruzada el poder predictivo del modelo.

En la Figura 4, se presentan los resultados de ejecutar el algoritmo MPPCA sobre los dos conjuntos de datos para el primer componente principal. Varias conclusiones pueden extraerse. En primer lugar, parece conveniente dividir el modelo del NOC del proceso nylon 6'6 en varios submodelos, ya que la mejora con la separación es considerable. El algoritmo llega a quintuplicar la varianza explicada por el modelo completo para una variable latente (de un 7,4% a un 37,2%). Este resultado es bastante sorprendente, más aún si tenemos en cuenta que la subdivisión únicamente detecta 4 submodelos. En la Figura 4(a) se observa que a partir de un umbral de un 14% el algoritmo deja de detectar subdivisión, y por tanto mejora. En el

caso del otro proceso la subdivisión no parece justificada, ya que las únicas divisiones encontradas (para un umbral de un 1%) no reflejan una mejora considerable (de un 8% a un 10%).

También se observa que se obtienen mejores resultados al reducir el tamaño mínimo de submodelo y al reducir el umbral. Este resultado es ciertamente lógico, debido al algoritmo utilizado, que realiza una búsqueda exhaustiva en cada recursión. Si una ejecución del algoritmo obtiene un resultado, otra que use un tamaño de submodelo menor obtendrá, al menos, el mismo resultado. Sin embargo, es posible que mejore, ya que dispone de un conjunto mayor de instantes de medición donde subdividir. Para un tamaño de submodelo de 1/8 se obtiene una mejora de un 11,8% de varianza

explicada (con 4 submodelos) con respecto a un tamaño de  $1/3$  (con 2 submodelos).

Por otro lado, si una ejecución del algoritmo obtiene un resultado, otra que use un umbral menor obtendrá, al menos, el mismo resultado. Sin embargo, también es posible que mejore, ya que dará por buenas algunas divisiones que el anterior no, aumentando la profundidad de la recursividad y aumentando la varianza explicada.

Ambas tendencias no son estrictas, ya que mientras que la medida de bondad final se obtiene por validación cruzada, la medida de bondad de una separación en el algoritmo no (ver sección 3.1). El hecho de que se haya mantenido el carácter monótono de la medida final, indica el buen ajuste que la medida de bondad de la separación no validada hace sobre la calculada con validación cruzada.

Analícemos ahora qué ocurre al aumentar el número de componentes principales. Para realizar el análisis se ha seleccionado un umbral del 7%. El umbral escogido no encuentra división alguna en el proceso de serigrafía, para ninguno de los valores de los otros parámetros. En la Figura 5 se presenta el resultado sobre los datos del proceso nylon 6'6. Nótese que la Figura 5(a) presenta la mejora en varianza explicada al utilizar MPPCA en lugar de MPCA, y no la varianza explicada total obtenida. Puede observarse el buen comportamiento del algoritmo al aumentar el número de componentes, manteniendo una mejora considerable y una cantidad de subfases prácticamente constante. La respuesta del algoritmo es muy positiva, ya que los modelos lineales con muchos componentes enmascaran la no linealidad, ofreciendo alta varianza explicada pero baja capacidad discriminante.

En consecuencia, los resultados permiten concluir que el algoritmo permite discriminar los conjuntos de datos lineales y no lineales a través del umbral de mejora. En el caso en que el conjunto de datos sea no lineal, el algoritmo genera un modelo multifase con una buena aproximación a los datos.

## 6. GENERACIÓN DE LAS TABLAS DE MONITORIZACIÓN

Una vez se ha obtenido el modelo estadístico del proceso, se utilizan un conjunto medidas para generar representaciones gráficas que simplifiquen las tareas de monitorización. Típicamente se utilizan la D-estadística (D-statistic) y la Q-estadística (Q-statistic) para la detección fallos, y la contribución de las variables como apoyo al diagnóstico (Nomikos and MacGregor, 1995). Las dos primeras medidas permiten generar unos límites de control a cierto nivel de confianza, de

manera que si éstos son traspasados por un lote se produce un aviso o detección de fallo, según el procedimiento empleado. Estos límites se sustentan en la suposición de que tanto las puntuaciones como los residuos siguen una distribución multinormal. Una vez un fallo ha sido notificado, las tablas de contribución permiten detectar las variables relacionadas con éste. Otra gráfica utilizada es la tabla de puntuaciones (t-plots), que permite observar la distribución de las puntuaciones en el espacio de las variables latentes.

Para la monitorización de procesos por lotes, la división en múltiples modelos encierra principalmente un problema: a mayor número de modelos, mayor número de tablas de control que hacen que la misión del operario sea más ardua. Para simplificar esta misión, Ündey y Çinar (Ündey and Çinar, 2002) proponen una estructura jerárquica basada en la construcción de una supermatriz de consenso. Una vez que se ha detectado un fallo, esta matriz permite reconocer qué fase ha contribuido en mayor medida. Así, el operario sabe qué tablas de control observar. Los autores también afirman que este paso se podría automatizar.

Esta solución no reduce el número de tablas, simplemente establece un procedimiento para indicar a qué tabla mirar. Desde este punto de vista, la propuesta no es demasiado útil. La separación en submodelos es temporal, por tanto, el resultado del comportamiento de un lote en una fase estará disponible al finalizar la propia fase. El operario únicamente tiene que observar las tablas de monitorización de una fase tras su finalización, de forma que no necesita que una matriz de consenso le indique dónde mirar. Además, este proceso se puede automatizar simplemente mostrando por pantalla las tablas de monitorización de la fase que haya acabado de terminar.

La propuesta defendida en este artículo es la siguiente: las puntuaciones obtenidas por el NOC en cada submodelo son escaladas según su matriz de varianza. Así, los límites de monitorización de todas las subfases se homogeneizan, de manera que las puntuaciones que un lote genera en distintas subfases pueden ser monitorizadas en una única tabla. Lo mismo ocurre con todas las estadísticas de un mismo tipo. El resultado es que el número de tablas es el mismo, la complejidad de las mismas también, pero en vez de observar un valor por cada lote, se observan tantos como subfases.

Para mostrar la generación del conjunto de tablas de monitorización a partir del modelo multifase y comprobar su rendimiento, se ha partido del conjunto de datos del proceso nylon 6'6. Se construirá un modelo con pocas subfases y dos componentes principales, para que los resultados sean de fácil interpretación. Utilizando la división obtenida para un tamaño de submodelo de  $1/3$ , se detectan



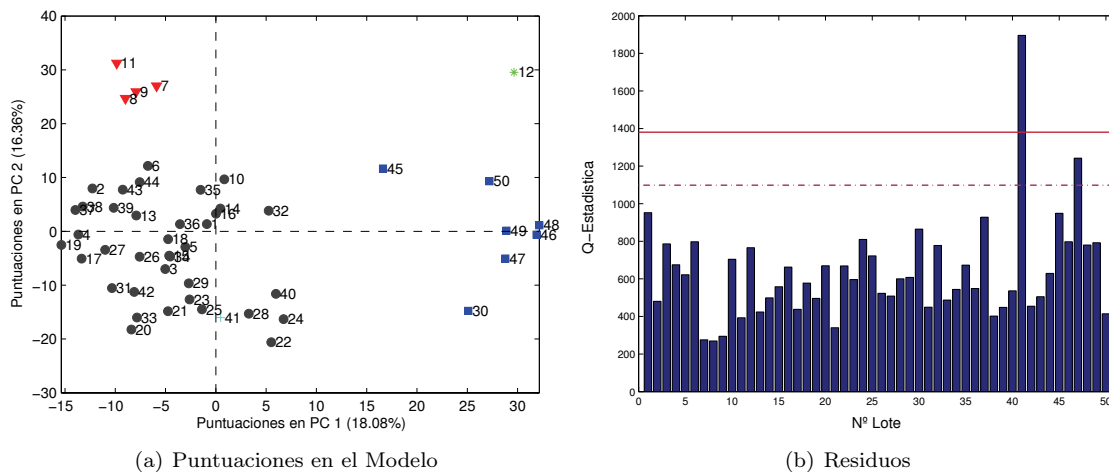


Figura 6. Agrupaciones en el NOC del proceso nylon 6'6, modelo completo.

dos subfases (de  $k = 1$  a  $k = 68$ , y de  $k = 69$  a  $k = 116$ ), produciendo una mejora absoluta del 12,3% de varianza explicada (del 26,5% a un 38,8%).

El conjunto de datos utilizado es el del Reactor A estudiado en (Kosanovich *et al.*, 1996). Sin embargo, los lotes no han sido alineados a la misma longitud, ni la subdivisión se realiza en el mismo punto, por lo que los resultados gráficos aquí obtenidos no serán idénticos. Una aproximación de la subdivisión propuesta en el artículo citado obtiene una mejora absoluta del 8,5%, casi un tercio menor a la obtenida aquí. Estos datos son, igualmente, los utilizados en (Nomikos and MacGregor, 1995). Sin embargo, en este artículo el orden de los lotes ha sido alterado para facilitar la comprensión del lector, lo que hace difícil la comparación.

Tras la detección de *outliers* en la siguiente sección, se introduce cómo se generan cada una de las tablas de detección de fallos comúnmente utilizadas y cómo han de ser modificadas para su uso a partir de un modelo multifase.

### 6.1 Detección de outliers

El modelado estadístico es sensible a los *outliers*. Por este motivo, es necesaria su detección y eliminación antes de generar las tablas de monitorización. El proceso de eliminación de *outliers* es un proceso iterativo en el que a través de la observación de la distribución de las muestras, tanto en el sub-espacio del modelo como en los residuos, se detectan y eliminan aquellas que no siguen el patrón general. Con el conjunto de datos resultante se procede al cálculo de los límites estadísticos que delimiten las condiciones de operación normal.

La Figura 6(a) representa las puntuaciones de los lotes en el espacio de los dos primeros componentes principales del modelo PCA. En la figura, podemos destacar un cluster principal formado por círculos (esquina inferior izquierda) y cuatro tipos de *outliers*. El primer tipo aparece a la derecha y está dibujado con cuadrados. El segundo aparece en la esquina superior izquierda, dibujado como triángulos invertidos. El tercer tipo está únicamente formado por el lote 12, que aparece ubicado en la esquina superior derecha. Por último, el lote 41 es un *outlier* detectable por los residuos (Figura 6(b), ver 13 para el cálculo de la Q-estadística).

Tras hacer la división en dos submodelos, los mismos conjuntos de *outliers* pueden ser observados. El primer grupo antes comentado (lotes: 30 y del 45 al 50), conjuntamente con el lote 12, aparecen claramente distanciados del cluster principal en la primera subfase, Figura 7(a). Este mismo resultado fue presentado por (Kosanovich *et al.*, 1996). El efecto de los *outliers* en el modelo PCA puede ser observado en esta figura. El conjunto de lotes normales presenta una clara tendencia a situarse a lo largo de una línea. Sin embargo, la dirección del primer componente principal calculada por el modelo se ve afectada por el grupo de outliers, haciendo necesaria además la adición de un segundo componente al modelo.

La segunda subfase es donde el segundo grupo de *outliers* (lotes: 7, 8, 9 y 11) se destaca del resto (Figura 8(a)). De nuevo, el lote 12 y los lotes 46 y 47, en distintas direcciones, vuelven a aparecer distanciados. El lote 41 también es identificado como outlier en los residuos de ambas subfases (Figuras 7(b) y 8(b)).

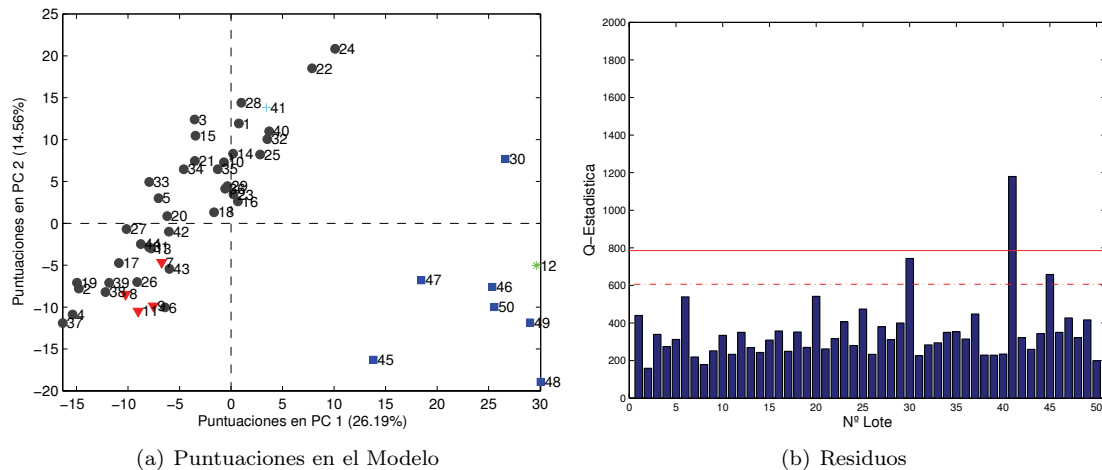


Figura 7. Agrupaciones en el NOC del proceso nylon 6'6, submodelo 1.

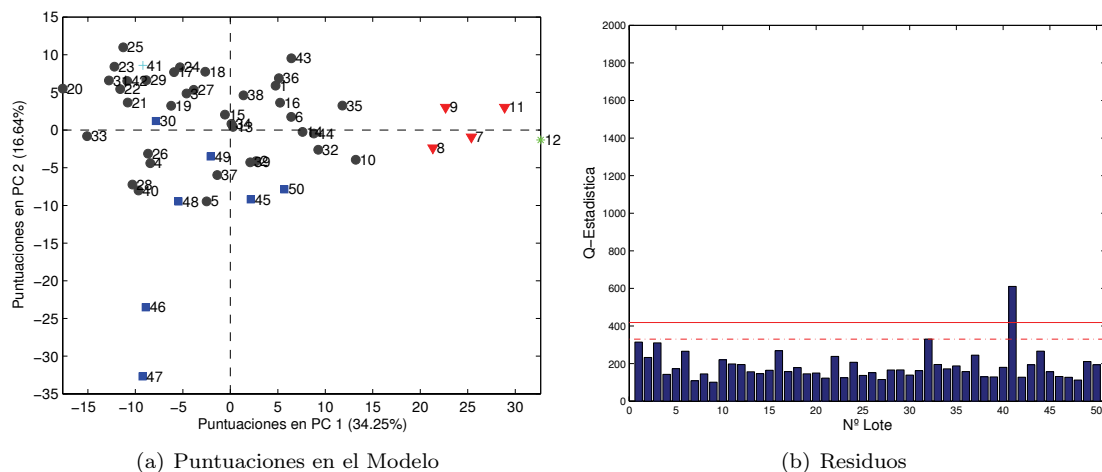


Figura 8. Agrupaciones en el NOC del proceso nylon 6'6, submodelo 2.

Una de las características más destacadas de las técnicas de análisis por variables latentes, como PCA, es su capacidad para ofrecer información de diagnóstico sobre un fallo o bien un *outlier*. Por ejemplo, en la primera subfase (Figura 7(a)) sabemos que el lote 12 se distancia del NOC principalmente por su valor en el componente principal 1 (nótese que su valor en el segundo componente está dentro de la normalidad). A partir de la contribución de las variables a su valor en este componente, se detecta que las variables 6, 8 y 9 son las principales implicadas. A partir de la contribución de las medidas temporales, se detecta que el lote se separó del comportamiento normal principalmente en el intervalo [18, 25].

Todos estos resultados confirman que con el modelo multifase pueden detectarse los mismos *outliers*

que con el modelo completo. Además, el modelo multifase proporciona información añadida: en qué subfase un lote se ha desviado del comportamiento normal.

### 6.2 Tablas de Componentes Principales

Cuando las puntuaciones obtenidas por un modelo PCA son representadas en el espacio de los componentes principales, es posible establecer unos límites de confianza con la premisa de que la distribución de las puntuaciones a lo largo de un componente principal es normal. Estos límites, para cierto valor de confianza, toman forma hiperelipsoidal en el hiperplano formado por el subespacio PCA. Los límites para un  $\alpha\%$  de confianza para

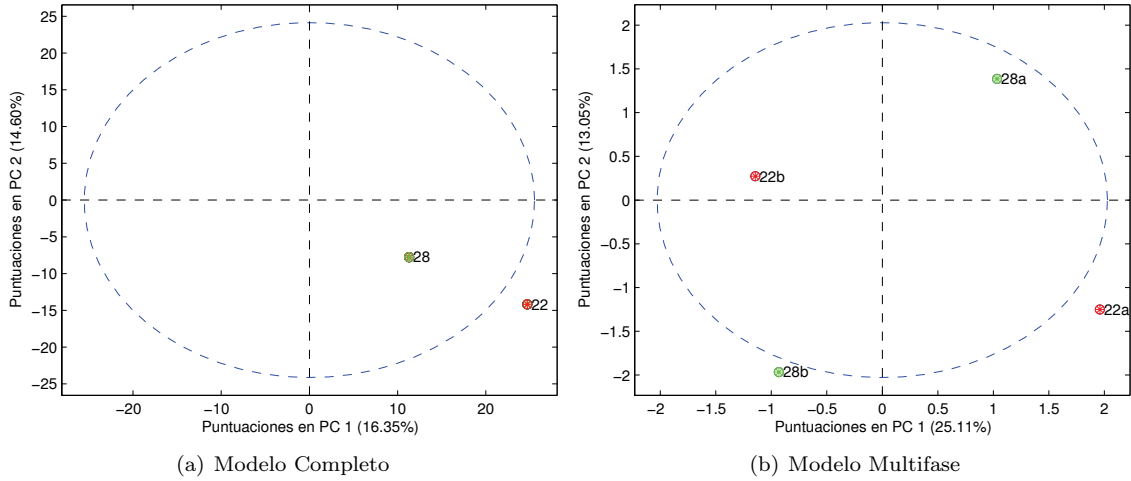


Figura 9. Tabla de Componentes Principales.

una nueva puntuación en un cierto instante, asumiendo normalidad e independencia, se calculan según (5).

$$\pm t_{n-1, \alpha/2} s_{ref} (1 + 1/n)^{1/2} \quad (5)$$

donde  $\pm t_{n-1, \alpha/2}$  es el intervalo centrado en el origen de una distribución t-student de  $n-1$  grados de libertad que encierra el  $\alpha\%$  de las muestras de esta distribución,  $s_{ref}$  es la desviación estándar estimada y  $n$  el número de observaciones, ambos medidos para la puntuación en el correspondiente instante.

Para que las puntuaciones obtenidas a partir de distintos modelos PCA (por ejemplo, que representen a distintas subfases) puedan ser monitorizadas con la misma gráfica, es necesario que los límites establezcan la misma región de normalidad. La propuesta defendida en este artículo es modificar la matriz de cargas según la ecuación (6).

$$\mathbf{P}_m = \mathbf{P}\mathbf{S}_d^{-1} \quad (6)$$

donde  $\mathbf{S}_d^{-1}$  es una matriz diagonal con los valores inversos de las desviaciones típicas de la distribución de puntuaciones en las direcciones de los componentes principales. Tras realizar esta modificación en la matriz de cargas, las puntuaciones del NOC han de ser recalculadas (ecuación 7) o bien multiplicadas a su vez por  $\mathbf{S}_d^{-1}$  (ecuación 8).

$$\mathbf{T}_m = \mathbf{X}\mathbf{P}_m \quad (7)$$

$$\mathbf{T}_m = \mathbf{T}\mathbf{S}_d^{-1} \quad (8)$$

Para recuperar los datos a partir del modelo, es necesario incorporar la matriz de varianzas de

las puntuaciones en las direcciones de los componentes principales (ver ecuación 9).

$$\mathbf{X} = \mathbf{T}_m \mathbf{S} \mathbf{P}_m \quad (9)$$

siendo  $\mathbf{S} = \mathbf{1}/(\mathbf{S}_d^{-1})^2$ .

El resultado es que la forma de la región de normalidad de las puntuaciones es ahora hiperesférica, y que las distribuciones de probabilidad en cada componente principal son de varianza la unidad. Las puntuaciones generadas por distintos modelos PCA que hayan sido modificados de esta manera pueden ser monitorizados a través de la misma tabla.

Como se observa en la Figura 9, la información conseguida a través del modelo multifase es mucho más completa, ya que permite detectar en qué subfase el comportamiento del lote ha variado notablemente con respecto al resto. En la figura 9(a), el lote 22 aparece fuera del límite de confianza, fijado al 95%. Sin embargo, no podríamos decir si se ha separado del resto de lotes en la primera o segunda subfase, o durante todo el lote. Con el modelo multifase se detecta que el comportamiento más alejado de la normalidad ocurre al inicio. En cuanto al lote 28, el modelo completo ni siquiera detecta que en la segunda subfase se sale del límite de confianza. Esto refleja que el modelo multifase posee un mayor poder de discriminación.

Si añadimos a la anterior ventaja el hecho de que no hace falta esperar a que el lote haya sido completado para disponer de datos que permitan monitorizarlo, esta variante de la monitorización ofrece una alternativa muy atractiva.

Nótese que la dirección en la que se sitúan las puntuaciones con respecto al centro no es significa-

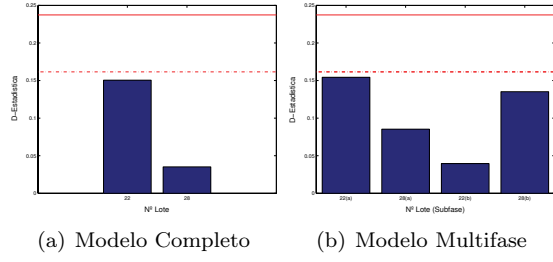


Figura 10. Tabla de D-Estadística.

tiva, ya que cada puntuación se obtiene tomando sistemas de referencia rotados.

### 6.3 D-estadística

La ecuación general para obtener la D-estadística es (10):

$$D = (\mathbf{t} - \bar{\mathbf{t}})' \mathbf{S}^{-1} (\mathbf{t} - \bar{\mathbf{t}}) \frac{I}{(I-1)^2} \quad (10)$$

donde  $\bar{\mathbf{t}}(R \times 1)$  es la puntuación media de la matriz  $\mathbf{T}$  calculada para el NOC con MPCA,  $R$  el número de componentes principales,  $\mathbf{S}(R \times R)$  es la matriz de covarianza de  $\mathbf{T}$  e  $I$  es el número de observaciones del NOC. Nótese que, mientras para Nomikos y Macgregor (Nomikos and MacGregor, 1994)  $\bar{\mathbf{t}}$  es igual a 0 y no aparece en la ecuación, en otras propuestas sí ha de ser considerada (van Sprang *et al.*, 2002).

Para obtener la D-estadística del modelo multifase, siguiendo la propuesta aquí defendida, partiremos de la matriz  $\mathbf{T}_m$ . El cálculo de la D-estadística a partir de las puntuaciones modificadas es directo, eliminando de la ecuación (10) tanto la puntuación media ( $\bar{\mathbf{t}}$ ) como la inversa de la matriz de covarianza ( $\mathbf{S}^{-1}$ ). La primera no es necesaria al utilizar el preprocesamiento propuesto en (Nomikos and MacGregor, 1994). La segunda tampoco lo es al haber modificado la matriz de las puntuaciones. El resultado queda expresado en la ecuación (11):

$$D = \frac{\mathbf{t}_m' \mathbf{t}_m I}{(I-1)^2} \quad (11)$$

Siguiendo el ejemplo anterior, en la figura 10 se observan las gráficas de cada modelo para los dos lotes seleccionados. Las líneas rojas continua y discontinua representan los límites al 99% y 95%, respectivamente, de la D-Estadística, calculados a partir de la distribución *Beta* (ver ecuación 12).

$$D \sim B(R/2, (I-R-1)/2) \quad (12)$$

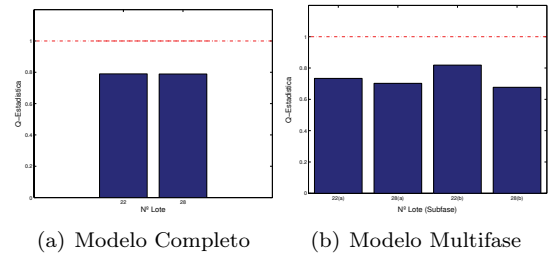


Figura 11. Tabla de Q-Estadística.

Como se puede observar, los resultados reflejan lo mismo que en la sección anterior, con la diferencia de que el límite impuesto en la D-Estadística es algo menos restrictivo que el obtenido con la distribución t-Student.

### 6.4 Q-estadística

La Q-estadística se define según la ecuación 13.

$$Q = \sum_{j=1}^J \sum_{k=1}^K e(j, k)^2 \quad (13)$$

$$e(j, k) = x(j, k) - \hat{x}(j, k)$$

donde  $x(j, k)$  es el valor medido de la variable  $j$  en el instante de medición  $k$ , y  $\hat{x}(j, k)$  es su predicción utilizando algún modelo estadístico, como MPCA.

Para obtener una única gráfica de monitorización para la Q-estadística a partir del modelo multifase, los valores obtenidos son escalados por el límite de confianza seleccionado (Louwerse and Smilde, 2000). De esta manera, todos los valores dentro de los límites de confianza aparecerán en el intervalo  $[0, 1]$ . Por encima del 1 los residuos serán entendidos como anormales.

De nuevo, la ventaja es obtener información más detallada del proceso. En la Figura 11(a), ambos lotes presentan una estadística muy similar. Sin embargo, su comportamiento en las fases no es el mismo (Figura 11(b)). Ambos lotes se mantienen dentro del límite de normalidad al 95%.

## 7. CONCLUSIONES

El presente artículo propone un esquema de monitorización de procesos por lotes a partir de múltiples modelos MPCA. Este esquema incluye la definición de un algoritmo para la detección de segmentos lineales a lo largo de la duración del lote, el algoritmo Multi-Phase PCA, y la modificación de las tablas de monitorización tradicionales.

Los resultados obtenidos muestran como algunos procesos, tras preprocesar los datos, siguen presentando no linealidades, que hacen que técnicas como PCA obtengan modelos pobres. En estos casos, el modelado mejora ampliamente al dividir convenientemente el lote en subfases, donde la estructura de correlación entre variables se mantiene suficientemente constante.

El algoritmo presentado, MPPCA, obtiene la división del modelo del proceso a partir de dos parámetros de entrada suficientemente intuitivos: el tamaño mínimo de submodelo y la mejora mínima obtenida para aceptar una subdivisión. El número de componentes principales de cada modelo local es el mismo que el del modelo completo. Así un diseñador, sin conocimiento específico del proceso, podría obtener el modelo multifase haciéndose dos simples preguntas: ¿qué tamaño mínimo de conjunto muestras temporales estoy dispuesto a modelar? y ¿cuánta mejora me justifica una subdivisión?. También se ha sugerido un valor del 7% como umbral apropiado independientemente del proceso a tratar.

En la segunda parte del artículo, se propone un conjunto de modificaciones para adaptar las tablas de monitorización más utilizadas al modelo multifase. El resultado es que la detección de outliers es la misma que con el modelo completo, el número de tablas y la complejidad de las mismas también lo es, pero la información obtenida es mayor.

En definitiva, el esquema de monitorización propuesto ofrece mejores resultados que el tradicional para los procesos no lineales. En cuanto a los procesos lineales, el análisis se reduce al MPCA tradicional. De esta manera, la propuesta es una generalización de la monitorización de procesos por lotes a través de MPCA. El modelo multifase resultante es de gran valor para mejorar el conocimiento sobre el proceso.

Una línea de trabajo interesante desde el punto de vista de la ingeniería de control consiste en el desarrollo de metodologías que, a partir de los resultados obtenidos en cada lote mediante las técnicas descritas, implementen la optimización de algún criterio de calidad a lo largo del *eje de lotes*.

## Apéndice A. NON-LINEAR ITERATIVE PARTIAL LEAST SQUARES (NIPALS)

El algoritmo NIPALS (Wold and Lyttkens, 1969) se muestra a continuación:

Para cada componente (1...A)

Inicializa  $\mathbf{t}$  con una columna de  $\mathbf{X}$

Repite hasta convergencia de  $\mathbf{t}$

$$\mathbf{p} = \mathbf{X}^T \cdot \mathbf{t}$$

Normalizar a  $\|\mathbf{p}\| = 1$

$$\mathbf{t} = \mathbf{X} \cdot \mathbf{p}$$

fin

$$\mathbf{X} = \mathbf{X} - \mathbf{t} \cdot \mathbf{p}^T$$

## Apéndice B. AGRADECIMIENTOS

Trabajo parcialmente financiado a través del programa de becas de Formación de Profesorado Universitario (FPU), Secretaría de Estado de Educación y Universidades, España. Los autores quisieran agradecer a Kenneth S. Dahl (DuPont) y a Neal B. Gallagher (Eigenvector) su amabilidad al suministrar los conjuntos de datos utilizados en el artículo.

## REFERENCIAS

- Abonyi, J., B. Feil, S. Németh and P. Arva (2005). Modified gath-geva clustering for fuzzy segmentation of multivariate time-series. *Fuzzy Sets and Systems* **149**, 39–56.
- Camacho, J. and J. Picó (2006). Multi-phase principal component analysis for batch processes modelling. *Chemometrics and Intelligent Laboratory Systems* **81**(2), 127–136.
- Díez, J.L., J.L. Navarro and A. Sala (2004). Control por planificación de ganancia con modelos borrosos. *Revista Iberoamericana de Automática e Informática Industrial* **1**(1), 32–43.
- Dong, D. and T.J. McAvoy (1996). Batch tracking via non-linear principal component analysis. *American Institute of Chemical Engineering Journal* **42**(8), 2199–2208.
- Edgar, T.F. (2004). Control and operations: when does controllability equal profitability?. *Computers and Chemical Engineering* **29**(1), 41–49.
- Kassidas, A., J.F. MacGregor and P.A. Taylor (1998). Synchronization of batch trajectories using dynamic time warping. *AIChE Journal* **44**, 864–875.
- Kosanovich, K.A., K.S. Dahl and M.J. Piovoso (1996). Improved process understanding using multiway principal component analysis. *Engineering Chemical Research* **35**, 138–146.
- Kourti, T. (2002). Process analysis and abnormal situation detection: from theory to practice.

- IEEE Control Systems Magazine* **22**(5), 10–25.
- Louwerse, D.J. and A.K. Smilde (2000). Multivariate statistical process control of batch processes based on three-way models. *Chemical Engineering Science* **55**, 1225–1235.
- Lu, N., F. Gao and F. Wang (2004). Sub-pca modeling and on-line monitoring strategy for batch processes. *AIChE Journal* **50**(1), 255–259.
- Nomikos, P. and J.F. MacGregor (1994). Monitoring batch processes using multiway principal components analysis. *AIChE Journal* **40**(8), 1361–1375.
- Nomikos, P. and J.F. MacGregor (1995). Multivariate spc charts for monitoring batch processes. *Technometrics* **37**(1), 41–59.
- Rotem, Y., A. Wachs and D.R. Lewin (2000). Ethylene compressor monitoring using model-based pca. *AIChE Journal* **46**(9), 1825–1835.
- Singhal, A. and D.E. Seborg (2002). Pattern matching in historical batch data using pca. *IEEE Control Systems Magazine* **22**(5), 53–63.
- Trelea, I.C., G. Trystarm and F. Courtois (1997). Optimal constrained non-linear control of batch processes: Application to corn drying. *Journal of Food Engineering* **31**, 403–422.
- Ündey, C. and A. Çinar (2002). Statistical monitoring of multistage, multiphase batch processes. *IEEE Control System Magazine* **22**(5), 40–52.
- van Sprang, E.N.M., H. Ramaker, J.A. Westerhuis, S.P. Gurden and A.K. Smilde (2002). Critical evaluation of approaches for on-line batch process monitoring. *Chemical Engineering Science* **57**, 3979–3991.
- Westerhuis, J.A., T. Kourti and J.F. MacGregor (1999). Comparing alternative approaches for multivariate statistical analysis of batch process data. *Journal of Chemometrics* **13**, 397–413.
- Wise, B.M., N.B. Gallagher, S.W. Butler, D.D. White Jr. and G.G. Barna (1999). A comparison of principal component analysis, multiway principal component analysis, trilinear decomposition and parallel factor analysis for fault detection in a semiconductor etch process. *Journal of Chemometrics* **13**, 379–396.
- Wold, H. and E. Lyttkens (1969). Nonlinear iterative partial least squares (nipals) estimation procedures. In: *Bull. Intern. Statist. Inst. Proc., 37th session*. pp. 1–15.
- Wold, S. (1978). Cross-validatory estimation of the number of components in factor and principal components. *Technometrics* **20**(4), 397–405.
- Wold, S., P. Geladi, K. Estebensen and J. Ohman (1987). Multi-way principal components and pls analysis. *Journal of Chemometrics* **1**, 41–56.