



**Pagura, José Alberto**

*Instituto de Investigaciones Teóricas y Aplicadas, Escuela de Estadística.*

## **EL USO DE COMPONENTES PRINCIPALES EN LA MONITORIZACION Y DIAGNOSTICO DE PROCESOS INDUSTRIALES**

### **1) INTRODUCCIÓN**

El control estadístico de los procesos, en su forma tradicional, analiza un pequeño conjunto de variables de calidad de forma individual y monitoriza el proceso, realizando a partir de las señales de fallo que se detecten, las acciones correctivas pertinentes. El estudio de las variables en forma separada omite el hecho de que una anomalía en el proceso puede evidenciarse no solo por valores inusuales en cada una de las características que se estudian, sino también por distorsiones en las relaciones entre ellas. Estas variables pueden analizarse en forma conjunta lo que permite un mejoramiento en la detección de anomalías. El tratamiento puede enriquecerse con la aplicación de métodos que permiten encontrar las variables responsables de tales situaciones.

Por otra parte, en muchas industrias en la actualidad, puede disponerse de información no sólo de variables de calidad de los productos sino también de una importante cantidad de datos del proceso, posiblemente relacionados con las características de calidad. La disponibilidad de estos datos puede ser aprovechada con diferentes propósitos: controlar el sistema monitorizando las variables del proceso, predecir valores de calidad del producto a través de las variables del proceso, determinar las condiciones en las que debe funcionar un proceso para obtener mejoras en la calidad, etc.

La aplicación de métodos multivariados que permiten la reducción de las dimensiones, entre ellos el análisis en componentes principales, aparecen como herramientas adecuadas para el mayor aprovechamiento de la información disponible.

En este trabajo, se discute el uso de componentes principales para la determinación de variables responsables de una señal de salida de control y la monitorización y diagnóstico de procesos, mostrando sus posibilidades y ventajas para la mejora de la calidad.

### **2) EL CONTROL MULTIVARIANTE DE PROCESOS**

El objetivo del control estadístico de procesos es monitorizar el comportamiento de los mismos con el fin de detectar cualquier evento inusual que desvíe el proceso de su funcionamiento normal (bajo control estadístico), pudiendo comprometer tanto la calidad del producto producido, como la seguridad del proceso de fabricación. Una vez identificada la causa que provocó la anomalía, es posible corregirla teniendo como consecuencia una mejora del proceso y por ende, de la calidad de los productos.

Los procedimientos que se utilizan en la práctica corriente, se basan en gráficos de control para un pequeño número de variables, tratándose por lo general de características de calidad del producto final. Los divulgados gráficos de Shewhart, CUSUM o EWMA son las herramientas para llevar adelante esta clase de control.

Cuando la calidad del producto está definida por varias variables, éstas pueden ser estudiadas en forma simultánea ya que una anomalía en el proceso puede estar anunciada no sólo por la detección de una salida de control en los gráficos correspondientes a algunas de las variables, sino también por una anomalía en la relación existente entre las variables que se estudian. Si se observan  $M$  características de calidad, distribuídas como una normal multivariante con vector de esperanzas  $\mu$  y matriz de covariancias  $\Sigma$ , que abreviaremos con  $N_M(\mu, \Sigma)$ , es posible dar un tratamiento basado en las distribuciones  $\chi^2$  y  $T^2$  de Hotelling.

Considérese un vector  $\mathbf{z}$  de  $M$  componentes correspondiente a las variables que se estudian, distribuído según  $N_M(\mu, \Sigma)$  en la situación "bajo control". La estadística :

$$\chi^2_o = (\mathbf{z} - \mu)' \Sigma^{-1} (\mathbf{z} - \mu) ,$$

se distribuye como una chi-cuadrado con  $M$  grados de libertad,  $\chi^2_{(M)}$ . Con esta referencia, ante una nueva observación se puede probar la hipótesis de pertenencia a dicha población normal comparando el  $\chi^2_o$  observado, con uno fijado según el nivel de significación deseado. Esto significa que se puede entonces fijar un límite superior de control igual al valor que antiacumula  $\alpha$  en la distribución chi cuadrado con  $M$  grados de libertad,  $\chi^2_{\alpha, M}$ , fijando  $\alpha$  en 0,01 o 0,05 y luego graficar las observaciones obtenidas en los diferentes momentos en los que se realizan las mediciones para el control, entendiéndose como señal de fuera de control a un valor de  $\chi^2_o$  superior al límite fijado.

Si se desconoce la matriz de covariancias  $\Sigma$ , puede estimarse a partir de una muestra calculando la matriz de covariancias muestral  $\mathbf{S}$ , y entonces considerar

$$T^2 = (\mathbf{z} - \mu)' \mathbf{S}^{-1} (\mathbf{z} - \mu)$$

que sigue una distribución  $T^2$  de Hotelling con  $m-1$  grados de libertad, donde  $m$  es el tamaño de la muestra con el que fuera calculada  $\mathbf{S}$ , actuando en forma análoga al caso anterior.

Una expresión útil para fijar el límite superior de control es la que relaciona a la distribución  $T^2$  con la  $F$  de Snedecor. Este enfoque permite, a partir de un conjunto de

$$T^2 = \frac{(N^2 - 1)M}{N(N - M)} F_{M, N - M}$$

variables de calidad del producto, detectar algún evento inusual en el proceso que afecte la calidad del producto.

Una forma de complementar la propuesta tradicional, que será tratada en el apartado siguiente, es utilizar procedimientos que permitan identificar las variables que son responsables de la señal de falta de control.

Otro aspecto a tratar, está sugerido por el hecho de que en muchos de los procesos industriales actuales se puede disponer de una importante cantidad de variables del proceso que, seguramente, están relacionadas con las variables de calidad del producto. Esta información se encuentra disponible en momentos anteriores a aquellos en los que pueden obtenerse las variables de calidad, y conociendo las condiciones operativas normales podrían monitorizar en forma anticipada al hallazgo de los fallos en las variables de calidad. Para ello, se puede establecer en las condiciones normales de operación, el comportamiento probabilístico de las mismas y establecer límites de control. Un valor observado superior al fijado como límite de control indicará una anomalía. En caso de conocer el valor poblacional de la matriz de covariancias, el límite se fijará a partir de la distribución  $\chi^2$ .

### 3) VARIABLES RESPONSABLES DE LAS SEÑALES DE SALIDAS DE CONTROL

Al encontrar una señal de salida de control, ya sea monitorizando las variables del proceso, o las variables de calidad, debe decidirse el curso de acción a seguir, para lo cual deben conocerse los factores que provocaron dicha anomalía. Kourti y Mac Gregor (1996) analizan diferentes propuestas para detectar qué variables son las responsables de una señal de salida de control y enuncian algunos métodos que serán considerados en este trabajo.

Llamando con  $\mathbf{z}$  al vector de observaciones de calidad o de observaciones del proceso, uno de los caminos posibles a seguir es calcular el error normalizado para cada variable de la observación  $i$ :  $(z_{ij}-\mu_j)/\sigma_j$ ,  $i=1\dots N$ , graficar sus valores para cada variable  $j$ , con  $j=1\dots M$  si se trata de variables de calidad o  $j=1\dots K$  si corresponden al proceso. Las variables que tengan grandes errores normalizados deberán ser estudiadas. Este análisis sobre los errores normalizados sirve para evidenciar salidas de control sólo por valores inusuales, dejando de lado las vinculadas a la correlación entre las variables. Esto significa, que puede haber observaciones que tengan los valores de las variables, en forma separada, dentro de los límites de control correspondientes a los diagramas univariados, pero estos valores pueden estar desproporcionados.

Una opción que mejora la propuesta anterior es la construcción de un gráfico de los llamados "scores" normalizados.

Los "scores" normalizados surgen de la posibilidad de descomponer la  $T^2$  utilizando los autovalores  $\lambda_j$ ,  $j=1\dots J$  de la matriz de covariancias  $S$  y combinaciones lineales de las variables  $\mathbf{z}$  a partir del siguiente razonamiento

$$\text{Sea } T^2 = (\mathbf{z}-\boldsymbol{\mu})' \mathbf{S}^{-1} (\mathbf{z}-\boldsymbol{\mu})$$

La matriz  $S$  puede diagonalizarse a partir de sus vectores propios, pudiéndose escribir como:

$\mathbf{S} = \mathbf{P} \mathbf{D}_\lambda \mathbf{P}'$ , siendo  $\mathbf{P}$  la matriz cuyas columnas son vectores propios normalizados de  $\mathbf{S}$ , y  $\mathbf{D}_\lambda$  matriz diagonal cuyos elementos  $d_{ii} = \lambda_i$  valor propio de  $\mathbf{P}$ . La matriz  $\mathbf{P}$  verifica:  $\mathbf{P}^{-1} = \mathbf{P}'$ .

$$\text{Entonces, } \mathbf{S}^{-1} = \mathbf{P} \mathbf{D}_\lambda^{-1} \mathbf{P}'$$

De esta forma,  $T^2$  se puede expresar como:

$$T^2 = (\mathbf{z}-\boldsymbol{\mu})' \mathbf{P} \mathbf{D}_\lambda^{-1} \mathbf{P}' (\mathbf{z}-\boldsymbol{\mu})$$

Llamando scores  $t_a$  a  $(\mathbf{z}-\boldsymbol{\mu})' \mathbf{P}$ , puede escribirse:

$$T^2 = \sum_{j=1}^J \frac{t_j^2}{\lambda_j} = \sum_{j=1}^J \frac{t_j^2}{s_{tj}^2}$$

Esta descomposición es llamada "descomposición en componentes principales". La estadística  $T^2$  queda expresado como una suma de razones donde los denominadores  $\lambda_j$ ,  $i=1\dots J$  son todos los autovalores no nulos de la matriz de covariancias  $\mathbf{S}$  y los numeradores  $t_i$ , llamados, scores. Se obtienen como combinaciones lineales de las  $\mathbf{z}$  con coeficientes iguales a las componentes de los autovectores normalizados correspondientes a cada autovalor ( $t_i = \mathbf{p}'_i \mathbf{z}$  donde  $\mathbf{p}'_i$  es el vector propio asociado a  $\lambda_i$ ). El número máximo de valores y vectores propios no nulos es el mínimo entre  $N-1$  y  $J$ .

Cada score  $t_a$  puede también expresarse como sigue:

$$t_j = p'_j(z - \mu) = \sum_{k=1}^J p_{j,k} (z_k - \mu_k)$$

Estos scores están distribuidos normalmente, por ser combinaciones lineales de variables normales con matriz de covariancias diagonal  $D_\lambda$ , es decir  $D_\lambda = \{d_{ij}\}$  con  $d_{ij} = \delta_{ij}\lambda_j$ ,  $j=1 \dots M$  y  $\delta_{ij} = 1$  si  $i=j$ , y es 0 en otro caso.

La estadística  $T^2$  obtenido a partir de los valores de  $z$  es igual a la que se obtiene con los scores de las componentes principales tomando todos los términos de la suma.

Los scores normalizados, es decir divididos por su desviación típica,  $\sqrt{\lambda_a}$ , se representan en un gráfico de barras. Aquéllos con valores altos son tenidos en cuenta para el estudio de causas asignables. Una forma de establecer un límite indicativo para decidir qué valor es alto, es fijarlos según el criterio de Bonferroni, para mantener el nivel global de significación fijado, es decir contrastarlos con el valor de tabla que corresponde al nivel  $\alpha/M$ , donde  $M$  es el número de scores.

Dado que los scores son combinaciones lineales de las variables originales  $z$ , cuando se pretenda interpretarlos con el fin de determinar que variables son responsables de la señal de anomalía, pueden presentarse problemas. Se propone entonces graficar, para aquellas observaciones con scores altos, las contribuciones de las variables en la composición de los mismos.

La **contribución de una variable  $z_j$  al score  $t_j$** , es  $p_{jj}(z_j - \mu_j)$ . Luego, las contribuciones grandes en valor absoluto y que llevan el mismo signo que el score considerado corresponderán a las variables que afectan directamente a la  $T^2$ , que permitió detectar la anomalía.

Es frecuente encontrar más de un score con valor absoluto alto, por lo que se deben estudiar las contribuciones de cada variable en forma separada. Para simplificar esta tarea y disponer de una herramienta para la rápida identificación, se propone el cálculo de las llamadas **contribuciones totales** de cada variable a todos los scores con valores altos, pero teniendo en cuenta solo aquellas variables con contribuciones de igual signo que el score obtenido. El cálculo se puede esquematizar en los siguientes pasos:

Para los  $K^*$  scores altos:

$$\text{cont}_{j',j} = \frac{t_{j'}}{s_{j'}} p_{j',j} (z_j - \mu_j)$$

Calcular la contribución de una variable  $z_j$  al score normalizado  $(t_{j'}/s_{j'})^2$

$\text{Cont}_{j'j}$  tomará el valor calculado si tiene el mismo signo que el score  $t_{j'}$  y será igual a 0 en caso contrario

Luego, calcular la contribución de la variable  $z_j$  como:

$$\text{CONT}_j = \sum_{j'=1}^J \text{Cont}_{j'j}$$

Esta última cantidad es la que permite identificar aquellas variables con altas contribuciones a la anomalía.

La utilización de las estadísticas  $\chi^2$  o de  $T^2$  puede presentar problemas en los casos de tener muchas variables y/o alta colinealidad.  $\Sigma$  o  $S$  pueden no ser invertibles. Por otra parte, la alta colinealidad provocará una mayor sensibilidad en las direcciones de los menores valores propios. Si se presentan datos faltantes también se dificulta el tratamiento. Una forma de tratar con el problema de la dimensionalidad puede ser subdividir el sistema en partes para las que tengan sentido los análisis por separado, pero esto no garantiza la solución a los problemas generados por la colinealidad.

Una solución aceptable para estas situaciones, la proporcionan algunos métodos de proyección como el análisis en componentes principales (PCA).

#### 4) ANALISIS EN COMPONENTES PRINCIPALES

Este método permite explicar la variabilidad existente en los datos de una matriz  $Z$ , de  $N$  observaciones y  $M$  variables, por medio de un número menor de variables latentes, combinación lineal de las primeras.

Desde el punto de vista geométrico, las filas de la matriz  $Z$  conforman una nube de puntos en el espacio  $M$  dimensional. Obtener las componentes principales es encontrar un subespacio de proyección sobre el cuál la nube dé la mínima deformación posible.

Si  $M < N - 1$ , pueden encontrarse hasta  $M$  componentes principales, es decir tantas como número de variables originales, pero el interés estará en seleccionar un número  $A$  mucho menor que  $M$  perdiendo lo menos posible de la variabilidad existente en el conjunto original de los datos.

La primera de estas variables latentes o primera componente principal  $t_1$  se obtiene buscando la combinación lineal de las variables originales  $t_1 = p_1 z$ , que tiene mayor variancia. Esta se expresa para el conjunto total de individuos como  $t_1 = Z p_1$  donde el vector  $p_1$  es el vector propio correspondiente al mayor valor propio  $\lambda_1$  de la matriz de covariancias de  $z$ ,  $\Sigma_z$ , y define la dirección de mayor variabilidad en el espacio de las  $M$  variables. Las componentes del vector  $p_1$  son denominadas cargas. El vector  $t_1$  de scores definido anteriormente, representa la proyección de cada observación sobre la dirección definida por  $p_1$ . Además, la variancia del score  $t_1$  es igual a  $\lambda_1$ , ya que:

$$\sigma_{t_1}^2 = p_1' E(z z') p_1 = p_1' \Sigma_z p_1 = p_1' \lambda_1 p_1 = \lambda_1$$

La segunda componente principal o variable latente es una combinación lineal de las variables originales  $t_2 = Z p_2$ , que se obtiene de forma tal que sea ortogonal a la primera, con la mayor variancia posible, inferior a la de  $t_1$ . Siguiendo los mismos criterios pueden encontrarse las restantes componentes principales.

Dada una observación  $i$ , los scores  $t_{i1}, t_{i2}, \dots, t_{iA}$ , ubican a dicha observación en el subespacio definido por los  $A$  primeros autovectores  $p_1, p_2, \dots, p_A$ , que son los asociados a los  $A$  mayores valores propios.

Si el cálculo de las componentes principales se realiza con los datos originales, las diferentes escalas en que están medidas las variables puede tener influencia decisiva en las direcciones que se obtienen. Por esta razón, previamente al cálculo de las componentes principales, los datos se centran y escalan restando a cada variable su media y dividiendo por su desviación estándar (aunque puede utilizarse otra forma de escalamiento).

Si se extraen todas las componentes principales,  $T = ZP$ , donde  $T$  es la matriz  $N \times M$ , cuyas columnas son los  $M$  vectores de scores y  $P$  es la matriz  $M \times M$  cuyas columnas son los

M vectores propios que constituye una base ortonormal. Postmultiplicando por  $\mathbf{P}'$  se tiene que:  $\mathbf{Z}=\mathbf{TP}'$ . Operando luego por bloques, la matriz  $\mathbf{Z}$  puede expresarse a partir de los scores  $\mathbf{t}$  y de los autovectores  $\mathbf{p}$ , de la siguiente forma:

$$\mathbf{Z} = \sum_{j=1}^J \mathbf{t}_j \mathbf{p}'_j = \sum_{j=1}^A \mathbf{t}_j \mathbf{p}'_j + \sum_{j=A+1}^J \mathbf{t}_j \mathbf{p}'_j = \hat{\mathbf{Z}} + \mathbf{E}$$

donde  $\hat{\mathbf{Z}}$  es la estimación de  $\mathbf{Z}$  utilizando solamente A componentes principales, asociadas a los A mayores valores propios, de las M posibles. Es decir, puede utilizarse un subconjunto A de componentes para reconstruir los valores originales con un cierto error.

Si bien el análisis en componentes principales es un método ampliamente divulgado con aplicación en diferentes campos, el uso de esta técnica para establecer un modelo de referencia no es tan habitual en la práctica. Este enfoque, sin embargo, es de gran utilidad para la monitorización y diagnóstico de procesos.

Otra forma equivalente de expresar el modelo, una vez que los datos han sido centrados y escalados se obtiene utilizando matrices en lugar de vectores:

$$\mathbf{Z} = \mathbf{1} \bar{\mathbf{z}}' + \mathbf{T} \mathbf{P}' + \mathbf{E}$$

El primer término es una matriz que contiene en cada columna la media de la variable correspondiente, el producto  $\mathbf{TP}'$  es la estructura del modelo, con  $\mathbf{T}$  matriz de scores,  $\mathbf{P}$  matriz de cargas, y  $\mathbf{E}=\{e_{ij}\}$  es la matriz de los residuales.

Cada fila de la matriz  $\mathbf{T}$  contiene las coordenadas de cada observación en el nuevo espacio definido por las componentes principales. Cada columna de  $\mathbf{P}$  contiene las cargas o pesos de cada una de las variables originales sobre cada variable latente. Las cargas definen la orientación del hiperplano obtenido, con respecto al espacio de las variables originales. Las cargas indican también la magnitud y sentido de la correlación de la componente considerada con cada una de las variables originales. En el caso en que las variables están tipificadas, la relación entre las cargas y el coeficiente de correlación de  $z_j$  y  $t_j$  está dado por:  $r(z_j, t_j) = p_{jj} \sqrt{\lambda_j}$

### **Cálculo de las componentes principales. El algoritmo NIPALS**

Los algoritmos para el cálculo de valores y vectores propios, de amplia divulgación, permiten obtener las componentes principales. Se pueden encontrar algoritmos que calculan todos los valores y vectores propios en forma simultánea y hay otros que lo hacen secuencialmente. Teniendo en cuenta que sólo se requerirá un número de componentes menor al número máximo que se puede calcular, dada la estructura de correlaciones que es normal encontrar entre las variables originales, conviene utilizar algún procedimiento que obtenga secuencialmente los vectores y valores propios desde el mayor hasta el más pequeño, indicando en qué momento no se desean más cálculos. Un algoritmo de uso habitual es el de Hotelling también conocido como "Power Method". Otro, que responde en forma adecuada, es el algoritmo NIPALS (Nonlinear Iterative Partial Least Squares). Cabe mencionar, que los algoritmos que obtienen todas las componentes en forma simultánea, son más precisos, pero pueden presentar inconvenientes si existe alta colinealidad ya que es posible que  $\mathbf{Z}'\mathbf{Z}$  esté mal condicionada. Los secuenciales son menos precisos pero funcionan bien en contextos de alta colinealidad y en presencia de datos faltantes.

Si bien en PCA cualquiera de los algoritmos desarrollados resultan adecuados, es de interés la utilización del NIPALS, teniendo en cuenta que, en el próximo método a estudiar, PLS, se utilizará también este algoritmo.

El algoritmo NIPALS puede sintetizarse en los siguientes pasos:

- 1) Elegir una columna de la matriz  $\mathbf{Z}$ , por ejemplo  $\mathbf{z}_j$  y hacer  $\mathbf{t}_{inicial} = \mathbf{z}_j$ .
- 2) Obtener un vector  $\mathbf{p}'_j = \mathbf{t}'_{inicial} \mathbf{Z} / (\mathbf{t}'_{inicial} \mathbf{t}_{inicial})$
- 3) Normalizar  $\mathbf{p}$
- 4) Calcular un nuevo  $\mathbf{t}_{final} = \mathbf{Z} \mathbf{p}_j$ . Las componentes de  $\mathbf{t}_{final}$  son las pendientes de la regresión de cada una de las filas de  $\mathbf{Z}$  sobre  $\mathbf{p}_j$ .
- 5) Comparar  $\mathbf{t}_{final}$  con  $\mathbf{t}_{inicial}$  y, estableciendo algún criterio de tolerancia, si resultan aproximadamente iguales se da por finalizada la búsqueda de la  $j$ -ésima componente. En caso contrario se debe repetir desde 2) asignando previamente el valor de  $\mathbf{t}_{final}$  a  $\mathbf{t}_{inicial}$  ( $\mathbf{t}_{inicial} \leftarrow \mathbf{t}_{final}$ )
- 6) Si se desea obtener otra componente, se extrae de la matriz  $\mathbf{Z}$  la parte explicada por la  $j$ -ésima componente ( $\mathbf{t}_j \mathbf{p}'_j$ ), es decir, se asigna a  $\mathbf{Z}$  el resultado de  $\mathbf{Z} - \mathbf{t}_j \mathbf{p}'_j$ , y se comienza nuevamente el proceso, es decir se vuelven a ejecutar los pasos 1 a 5.

En la convergencia se tiene que:

$$\mathbf{p}_j = \mathbf{Z}' \mathbf{t}_j / \mathbf{t}'_j \mathbf{t}_j = \mathbf{Z}' \mathbf{Z} \mathbf{p}_j / \mathbf{t}'_j \mathbf{t}_j ; \text{ luego, } \mathbf{Z}' \mathbf{Z} \mathbf{p}_j = \mathbf{t}'_j \mathbf{t}_j \mathbf{p}_j ,$$

luego  $\mathbf{p}_j$  converge al vector propio correspondiente al mayor valor propio de  $\mathbf{Z}' \mathbf{Z}$ . El valor propio será  $\mathbf{t}'_j \mathbf{t}_j$ .

La variancia del score  $\mathbf{t}_j$  encontrado, es el mayor autovalor de  $\mathbf{Z}' \mathbf{Z} / (N-1)$ , si se trata de una muestra, debiendo dividir entonces  $\mathbf{t}'_j \mathbf{t}_j$  por  $(N-1)$  para tener el valor propio buscado.

## HERRAMIENTAS DE DIAGNÓSTICO

PCA presenta interesantes herramientas de diagnóstico. Las mismas permiten detectar agrupamientos en las observaciones y en las variables así como tendencias, datos anómalos, etc.

Graficar los scores en los planos definidos por las variables latentes, permitirá encontrar agrupamientos y detectar datos con anomalías severas ("outliers" severos). Al representar gráficamente los scores junto con la elipse definida por la  $T^2$  de Hotelling, estos datos anómalos aparecen a una distancia del centro de la elipse mayor que cualquier punto de ella. Esta última comparación se propone a partir de la descomposición de la  $T^2$  en la suma de los scores. Además una observación con aporte  $t^2_{ij}$  elevado puede ser analizada observando a qué componente se debe ese aporte y adicionalmente, qué variables están relacionadas con tal componente.

Graficar las cargas permitirá visualizar las asociaciones entre las variables. La representación conjunta de observaciones y variables permitirá encontrar explicaciones sobre qué variables caracterizan los agrupamientos y tendencias hallados en las observaciones.

Es posible que, además de existir datos con anomalías severas, haya datos con anomalías moderadas. Una observación con estas características podrá tener coordenadas

interiores a la elipse definida por la  $T^2$  de Hotelling para un nivel de confianza dado, pero no estará bien ajustada por el modelo, es decir estará alejada del plano de ajuste de la nube. Esta situación estará reflejada en los residuales expresados en la matriz  $\mathbf{E}$  de la ecuación del modelo. Una observación en estas condiciones tendrá, en su correspondiente fila de la matriz  $\mathbf{E}$  valores absolutos grandes. Para la detección de "outliers" moderados, se puede utilizar la medida  $D_{\text{mod}Z_i}$  que representa la desviación típica residual del individuo  $i$  al modelo PCA ajustado  $\left(\sqrt{\frac{\text{SCR}}{gl}}\right)$  y que, como se observa, es proporcional a la raíz cuadrada de la distancia euclídea del individuo  $i$  a su proyección en el modelo.:

$$D_{\text{Mod}Z}_i = \sqrt{\frac{\sum_{j=1}^J e_{ij}^2}{J - A}}$$

donde  $e_{ij}$  es la diferencia entre la observación  $z_{ij}$  y su estimación por el modelo.

Si se define una distancia media de una observación al modelo como:

$$s_0 = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^A e_{ij}^2}{(N - A - 1)(J - A)}}$$

Puede calcularse la estadística  $(D_{\text{mod}Z_i}/s_0)^2$ , llamado también  $D_{\text{mod}Z_{\text{norm}}}$  (distancia al modelo normalizada) que sigue una distribución  $F_{J-A, (N-A-1)(J-A)}$ . De esa manera, comparando para una observación dada, el valor de  $D_{\text{mod}Z_{\text{norm}}}$  con la raíz cuadrada del percentil  $(1-\alpha)\%$  de una distribución  $F_{J-A, (N-A-1)(J-A)}$ , llamada distancia crítica normalizada, se puede probar si la distancia al modelo de una observación puede considerarse o no anómala. Si  $D_{\text{mod}Z_{\text{norm}}} > \sqrt{F_{J-A, (N-A-1)(J-A)}^\alpha}$  la observación  $i$ , es anómala moderada, es decir, no se ajusta bien al modelo PCA.

Es también importante diagnosticar qué variables están bien explicadas por el modelo. Para cada variable  $z_j$ , puede calcularse el % de variabilidad de esa variable, explicada por el modelo con  $A$  componentes,  $R^2_j$ ,

$$R^2_j = 1 - \frac{\text{SCR}_j}{\text{SCT}}$$

donde  $\text{SCR}_j = \sum_{i=1}^N e_{ij}^2$ , y  $\text{SCT} = \sum_{i=1}^N z_{ij}^2$  y  $e_{ij}$  son los residuos obtenidos al extraer  $A$  componentes.

Este valor puede variar entre 0, que corresponde al caso en el que el modelo no explique nada de la variable, y 1 para la situación contraria en la que toda la variabilidad de  $z_j$  está recogida por el modelo. Al aumentar el número de componentes principales, los valores de  $R^2_j$  irán aumentando hasta llegar a 1 cuando todas las variables latentes se incluyen en el modelo. Puede calcularse también el porcentaje de variabilidad total de  $\mathbf{Z}$  explicada por el modelo,  $R^2$ , con  $A$  componentes, como:



$$R^2 = 1 - \frac{SCR}{SCT}$$

$$\text{donde } SCR = \sum_{i=1}^N \sum_{j=1}^J e_{ij}^2 \text{ y } SCT = \sum_{i=1}^N \sum_{j=1}^J z_{ij}^2$$

Estos coeficientes  $R^2$  pueden ajustarse dividiendo cada suma de cuadrados por sus correspondientes grados de libertad obteniendo  $R^2_{aj}$  que mide el porcentaje de variancia explicada.

Cabe también analizar el número de componentes principales que debe incluirse en el modelo. La discusión se genera porque no es lo mismo considerar la capacidad del modelo para la representación simplificada de los datos reales o bondad de ajuste, que su capacidad de predecir nuevos datos. La primera de las características está medida por el valor  $R^2$  y a medida que se complica el modelo agregando variables latentes, su valor aumenta. Sin embargo, en la medida en que el modelo esté sobreparametrizado, su capacidad de predicción disminuirá. La capacidad de predicción puede ser evaluada mediante la técnica de validación cruzada. Esto es:

- 1) Dividir la muestra de N observaciones en L grupos de l unidades cada uno
- 2) Eliminar cada uno de los grupos por turno
- 3) Obtener el modelo con las N-l unidades restantes y luego predecir los valores de  $\mathbf{z}$  para las unidades eliminadas.

Para cada unidad i se podrá calcular el error cuadrático de predicción:  $\sum_{j=1}^J (z_{ij} - \hat{z}_{ij})^2$

donde  $\hat{z}_{ij}$  es el valor de la j-ésima variable para el i-ésimo individuo según la estimación provista por el modelo. La suma de cuadrados de los residuos de las predicciones para todas las observaciones se denomina PRESS (Predictive Residual Sum of Squares). A partir del PRESS se pueden construir diferentes indicadores de la bondad de predicción.

Si bien se cumple siempre que  $SCR(j) < SCR(j-1)$ , sin embargo,  $PRESS(j)$  será menor que  $SCR(j-1)$  sólo si la componente j-ésima mejora la capacidad predictiva del modelo, y en ese caso  $Q^2 > 0$ . En caso contrario  $PRESS(j) > SCR(j-1)$  y  $Q^2 < 0$ .

$$Q_j^2 = 1 - \frac{PRESS(j)}{SCR_{j-1}}$$

$SCR_{j-1}$  es la suma de cuadrados de los residuos del modelo con j-1 componentes y  $PRESS(j)$  es la suma de los cuadrados de los residuos obtenidos por validación cruzada en el modelo con j componentes.  $Q_j^2$  mide la variación de  $\mathbf{Z}$  predicha por la componente j-ésima.

$Q^2(\text{cum})_j = 1 - PRESS/SCT$  mide la bondad de predicción del modelo PCA con j componentes.  $Q^2$  indica la fracción de la variación total que puede ser predicha por una componente.  $Q^2$  no varía entre 0 y 1 como el caso de  $R^2$ ; el indicador de bondad de predicción crece desde 0 hasta un punto máximo donde comienza a decrecer aunque aumente el número de componentes, indicando en este caso un sobreajuste del modelo.

$Q^2$  puede ser calculado también para cada variable, lo que permite evaluar cuán bueno es el modelo para predecir los valores de cada una de las variables por separado. Un criterio práctico para decidir el número de variables latentes a utilizar, sugerido por Wold establece que si  $Q^2 > 0,5$  es buena, y  $Q^2 > 0,9$  es excelente. Además, la diferencia entre  $R^2$  y  $Q^2$  no debe ser muy grande.

### APLICACIÓN A LA MONITORIZACIÓN

Para la monitorización del proceso el procedimiento consiste en registrar, datos en condiciones normales de operación, y analizar los mismos decidiendo el número  $A$  de componentes principales a utilizar, obteniendo así la matriz de cargas  $\mathbf{P}$ . Ante una nueva observación  $\mathbf{z}_i$  se calculará su valor de  $T^2$

$$T_{i,A}^2 = \sum_{j=1}^A \frac{t_{i,j}^2}{S_{t_j}^2} = \sum_{j=1}^A \left[ \frac{\mathbf{p}'_j (\mathbf{z}_i - \boldsymbol{\mu})}{S_j} \right]^2$$

$$T_A^2 \sim \frac{A(N^2 - 1)}{N(N - A)} F_{A, N-A}$$

Ese valor se compara con el límite superior establecido, que corresponde al percentil  $(1 - \alpha) \%$  de una distribución  $F_{A, (N-A)}$  por la constante de la expresión anterior.

Debe también controlarse la distancia al modelo, ya que aunque la  $T^2$  calculada para una nueva observación, esté dentro de la región de aceptación, la misma puede estar mal ajustada por el modelo. Para ello, se puede calcular la distancia al modelo normalizada de la nueva observación y, compararla con distancia crítica normalizada. Para su aplicación práctica, se pueden construir gráficos de control de  $T^2$  y de  $(D_{modZ}/S_0)$ , cada uno con su límite superior de control. Ante una anomalía proceder al diagnóstico utilizando alguna de las herramientas propuestas.

### 5) COMENTARIOS FINALES

La gran cantidad de información que generan muchos de los procesos actuales, orienta hacia el desarrollo de métodos que permitan un mayor aprovechamiento de la misma. En primer lugar, los diagramas de control multivariados resultan una herramienta útil para la detección de salidas de control. En segundo lugar, para tomar las acciones correctivas adecuadas se deben encontrar las variables que provocaron dicha señal, lo cual no es tan evidente como en el control estadístico univariado. Esta tarea no es sencilla y se pueden encontrar en publicaciones recientes, diferentes propuestas, cuya implementación de manera inmediata, no parece fácil. No se cuenta aún, en el software estadístico de mayor divulgación la posibilidad de un fácil acceso a estas herramientas, sino que deben ser programadas.

Por otra parte, situaciones como el estudio de un gran número de variables, con alta colinealidad y datos faltantes, debilitan los instrumentos mencionados. La alternativa de la utilización del análisis de componentes principales, se constituye en un método apto para el tratamiento de datos con estas características, permitiendo el desarrollo de herramientas gráficas y medidas adecuadas, para la monitorización y diagnóstico de procesos. Nuevamente la mayor dificultad actual, es la posibilidad de una implementación rápida y



accesible, ya que la obtención de gráficos y medidas necesarias, no es inmediata en el software estadístico corriente, sino que deben desarrollarse los programas, encontrando solo en software más específicos como SIMCA<sup>1</sup> y PLS\_Toolbox<sup>2</sup> la posibilidad de una instrumentación sencilla en los procesos de producción.

### **Bibliografía**

- Fuchs, Camil Kenett, Ron S. "Multivariate Quality Control. Theory and Applications", Marcel Dekker Inc. , 1998
- Geladi, P. Y Kowalski, B. R. "Partial Least Squares Regression: A Tutorial". Analytica Chimica Acta 185, pág. 1-17, 1986
- Hodouin D, MacGregor J.F. Hou M. Franklin M. "Multivariate statistical analysis of mineral processing plant data". CIM Bulletin, Vol 85, N° 975, pág. 23-34, 1993
- Kourti T. and MacGregor J. F. "Multivariate SPC Methods for Process and Product Monitoring Journal of Quality Technology" Vol 28 N° 4, 1996
- Wold S. Multi and Megavariate Data Anlysis using Projection Methods
- SIMCA –P 8.0 User Guide and Tutorial UMETRICS, 1999

---

<sup>1</sup> SIMCA-P, UMETRICS AB, Box 7960, S-907 19 Umea, Sweden

<sup>2</sup> PLS\_Toolbox, Barry M. Wise-Neal B. Gallagher. Eigenvector Technologies, P.O. Box 483, 196 Hyacinth Avenue, Manson, WA, USA