



**UNIVERSIDAD DE CHILE
FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS
DEPARTAMENTO DE INGENIERIA ELÉCTRICA**

**DETECCIÓN DE ANOMALÍAS EN PROCESOS INDUSTRIALES USANDO
MODELOS BASADOS EN SIMILITUD**

**MEMORIA PARA OPTAR AL TÍTULO DE INGENIERO CIVIL
ELECTRICISTA**

ALEJANDRO SAMIR LEÓN OLIVARES

**SANTIAGO DE CHILE
ENERO 2012**



**UNIVERSIDAD DE CHILE
FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS
DEPARTAMENTO DE INGENIERIA ELÉCTRICA**

**DETECCIÓN DE ANOMALÍAS EN PROCESOS INDUSTRIALES USANDO
MODELOS BASADOS EN SIMILITUD**

**MEMORIA PARA OPTAR AL TÍTULO DE INGENIERO CIVIL
ELECTRICISTA**

ALEJANDRO SAMIR LEÓN OLIVARES

**PROFESOR GUÍA:
MARCOS ORCHARD CONCHA**

**MIEMBROS DE LA COMISIÓN
HECTOR AGUSTO ALEGRÍA
DORIS SÁEZ HUEICHAPAN**

**SANTIAGO DE CHILE
ENERO 2012**

RESUMEN DE LA MEMORIA PARA OPTAR
AL TÍTULO DE INGENIERO CIVIL
ELECTRICISTA.
POR: ALEJANDRO LEÓN O.
FECHA: ENERO 2012
PROF.GUÍA: DR. MARCOS ORCHARD C.

“DETECCIÓN DE ANOMALÍAS EN PROCESOS INDUSTRIALES USANDO MODELOS BASADOS EN SIMILITUD”

La detección de anomalías en procesos industriales es un tema de alto impacto que ha sido analizado y estudiado en diversas áreas de la ingeniería e investigación. La mayor parte de los métodos de detección actualmente disponibles posibilitan el estudio de las irregularidades encontradas en el historial de un proceso, ayudando a extraer información significativa (y a veces crítica) en una amplia variedad de aplicaciones, y convirtiéndose de este modo en parte fundamental e integral de esquemas de reducción de costos tanto humanos como económicos en la industria contemporánea.

El objetivo general de este trabajo es desarrollar e implementar un enfoque modular de detección de anomalías, aplicable a procesos industriales multivariados y fundado en el análisis de residuos generados a partir de modelos no paramétricos basados en similitud (*similarity-based modeling*, SBM). La herramienta consiste principalmente de un sistema de generación automática de modelos SBM, una metodología de estudio de eventos y un algoritmo estadístico de detección.

El trabajo realizado se enmarca en un proyecto de colaboración conjunta entre las empresas CONTAC, INGENIEROS LTDA. y ENDESA-CHILE. Gracias a esto, ha sido posible evaluar el sistema propuesto utilizando datos de operación correspondientes a una central termoeléctrica de ciclo combinado perteneciente a la segunda empresa ya mencionada.

Las comparaciones en cuanto al desempeño del sistema de modelación implementado permiten concluir que el algoritmo es capaz de generar una representación más apropiada del proceso, basado en el error obtenido con la técnica de modelación SBM, la cual es cercana al 25% del error obtenido con la técnica de modelación lineal en los parámetros.

Además, la metodología de estudio de eventos permite detectar correctamente las variables que no aportan al objetivo de detección de un evento en particular, así como también identifica las variables más significativas para lograr tal propósito, reduciendo el número de variables analizadas y con esto, los requerimientos de cómputo de operación en línea.

La validación de los resultados entregados por el método de detección de anomalías desarrollado, permite aseverar que la utilización de modelos no-paramétricos tipo SBM, en combinación con la metodología de estudio de eventos y el algoritmo estadístico de detección, resultan eficaces a la hora de generar alarmas y detectar las anomalías estudiadas.

AGRADECIMIENTOS

En primer lugar, agradezco a mis padres que forjaron en mí el espíritu de superación, compromiso, responsabilidad y lucha constante por alcanzar mis sueños. Así como también, agradezco a mis hermanos quienes me han acompañado durante gran parte de mi vida, gracias por su apoyo y paciencia.

A mi Profesor Guía, Dr. Marcos Orchard, por su admirable sabiduría, excelente disposición, su constante y fundamental apoyo en el desarrollo de este trabajo y el recibido durante mis últimos años de carrera.

A mis compañeros de proyecto, Sebastián Fuentealba y Hugo Navarrete, por toda la ayuda brindada durante el desarrollo del trabajo.

Al Grupo de Control Automático (GCA) de la Universidad de Chile y a todos mis amigos de Universidad, que de alguna u otra forma influyeron en mi formación como persona e ingeniero.

ÍNDICE DE CONTENIDOS

Capítulo 1. Introducción	1
1.1 Motivación.....	1
1.2 Alcance.....	2
1.3 Objetivos	2
1.4 Indicación sobre confidencialidad.....	3
1.5 Estructura general.....	3
Capítulo 2. Revisión bibliográfica	5
2.1 Técnicas de monitoreo de procesos: Conceptos básicos.....	5
2.2 Técnicas de detección de anomalías.....	8
2.3 Detección de anomalías basado en residuos.....	9
2.3.1 Métodos de generación de residuos	10
2.3.2 Generación de residuos basado en observadores	11
2.3.2.1 Modelos paramétricos	12
2.3.2.1.1 Modelos lineales-en-los-parámetros.....	13
2.3.2.1.2 Estimación de parámetros en modelos lineales-en-los-parámetros.....	16
2.3.2.2 Modelos no paramétricos.....	19
2.3.2.2.1 Modelos basados en similitud	20
2.3.3 Metodología de análisis de residuos	22
2.3.3.1 Análisis de componentes principales.....	24
2.3.3.2 Test estadístico de Hotelling.....	27
Capítulo 3. Implementación de herramienta de detección de anomalías usando modelos SBM... 30	
3.1 Introducción a las técnicas utilizadas	30
3.2 Algoritmo de generación automática de modelos SBM.....	34
3.2.1 Primer enfoque: Descripción general:.....	35
3.2.2 Segundo enfoque: Descripción detallada de funciones:.....	38

3.2.2.1	Bloques funcionales utilizados en identificación de procesos industriales usando modelos SBM.....	38
3.2.2.1.1	Pre-procesamiento (A).....	39
3.2.2.1.2	Limpieza de datos (B)	40
3.2.2.1.3	Análisis inicial (C).....	40
3.2.2.1.4	División conjunto inicial C1 (D)	41
3.2.2.1.5	SBM (E)	42
3.2.2.1.6	Similitud (F)	42
3.2.2.1.7	Vectores LD (G).....	42
3.2.2.1.8	Error de estimación (H).....	43
3.2.2.1.9	Filtro mediana móvil (I)	44
3.2.2.1.10	Análisis mediante Hotelling (J).....	44
3.2.2.1.11	Selección vía Hotelling (K).....	44
3.2.2.1.12	Ajuste final (L)	45
3.2.2.2	Bloque Principal: MODELO SBM.....	45
3.3	Metodología para estudio de eventos usando modelos SBM.....	50
3.3.1	Extracción de datos	51
3.3.2	Análisis inicial de variables	51
3.3.3	Creación de modelo	52
3.3.4	Re-ajuste/Validación.....	52
3.3.5	Prueba anomalía.....	52
3.3.6	Selección de variables.....	53
3.3.7	Modelo final.....	53
3.4	Detección de anomalías en procesos industriales usando modelos basados en similitud.....	53
Capítulo 4. Estudio de Evento – Resultados.....		56
4.1	Descripción del proceso industrial estudiado y de las bases de datos utilizadas.....	57
4.1.1	Central térmica de ciclo combinado.....	57

4.1.2 Bases de datos utilizadas.....	59
4.2 Resultados obtenidos en prueba de algoritmo de generación de modelos SBM y base de datos T1.....	62
4.3 Resultados obtenidos en la detección de eventos utilizando algoritmos de modelamiento basados en similitud y base de datos T2.....	80
Capítulo 5. Conclusiones	94
Bibliografía	96

ÍNDICE DE FIGURAS

Figura 2.1:	Esquema de ciclo diagnóstico de anomalías/fallas y recuperación (acomodación) del proceso.	8
Figura 2.2:	Proyección del vector de observación x , en los espacios score y de residuos, y el cálculo del vector de observación transformado.	27
Figura 2.3:	Ilustración grafica de la transformación del estadístico de Hotelling.	28
Figura 3.1:	Representación de observaciones transformadas por PCA (izquierda), selección de observaciones representativas (derecha).	32
Figura 3.2:	Ilustración de método de eliminación de variables basado en análisis de componentes PLS. Representación de variables en espacio transformado.	34
Figura 3.3:	Diagrama de flujo de algoritmo - Descripción general.	35
Figura 3.4:	Observaciones de la base de datos de entrenamiento graficados en el espacio transformado vía PCA y selección de aquellas observaciones representativas que constituyen C1.	41
Figura 3.5:	Subdivisión del conjunto inicial C1.	47
Figura 3.6:	Diagrama de Flujo – Descripción detallada de algoritmo de generación automática de modelos SBM.	49
Figura 3.7:	Diagrama de flujo de metodología para el estudio de eventos usando modelos SBM.	50
Figura 3.8:	Diagrama de flujo de herramienta de detección de anomalías usando modelos SBM.	54
Figura 4.1:	Esquema básico de Turbina a Gas/Diesel de una central térmica de ciclo combinado.	58
Figura 4.2:	Potencia activa T1 [MW].	60
Figura 4.3:	Velocidad turbina T1 [%].	60
Figura 4.4:	Referencia señal de combustible T1, DIESEL (en negro) y GAS (en verde).	61
Figura 4.5:	Potencia activa T2 [MW].	61
Figura 4.6:	Velocidad turbina T2.	61
Figura 4.7:	Referencia señal de combustible T2, DIESEL (en negro) y GAS (en verde).	62
Figura 4.8:	Observaciones proyectadas en la 1° y 2° componente principal de PCA.	63
Figura 4.9:	Observaciones clasificadas por modos de operación vía PCA.	64

Figura 4.10:	Comparación entre condiciones de operación identificables en a) y b), y la clasificación automática para cada observación c).	65
Figura 4.11:	Observaciones seleccionadas por análisis inicial.	66
Figura 4.12:	Resultado Preliminar obtenido con algoritmo de generación automática de modelos SBM.	67
Figura 4.13:	Aproximación a observación 5000 de la Figura 4.12.	68
Figura 4.14:	Resultados de modelo SBM inicial (4% de base de datos de entrenamiento). Gráfico c) en escala logarítmica.	69
Figura 4.15:	Resultados modelo SBM preliminar (10% de base de datos de entrenamiento). Gráfico c) en escala logarítmica.	71
Figura 4.16:	Resultado modelo SBM en 1° iteración. Gráfico c) en escala logarítmica.	71
Figura 4.17:	Resultado modelo SBM en 2° iteración. Gráfico c) en escala logarítmica.	72
Figura 4.18:	Resultado modelo SBM en 3° iteración. Gráfico c) en escala logarítmica.	72
Figura 4.19:	Resultado modelo SBM en 4° iteración. Gráfico c) en escala logarítmica.	73
Figura 4.20:	Resultado modelo SBM Final. Gráfico c) en escala logarítmica.	73
Figura 4.21:	EMC (Error Medio Cuadrático) de las observaciones estimadas vía modelo SBM final.	74
Figura 4.22:	Resultado modelo SBM final ajustado. Gráfico c) en escala logarítmica.	74
Figura 4.23:	EMC por observación de modelos SBM y LMV para base de datos de entrenamiento 1. Datos saturados superiormente (máximo 0.5).	78
Figura 4.24:	EMC por observación de modelos SBM y LMV para base de datos de validación 1. Datos saturados superiormente (máximo 0.4).	78
Figura 4.25:	EMC por observación de modelos SBM y LMV para base de datos de entrenamiento 2. Datos saturados superiormente (máximo 0.3).	79
Figura 4.26:	EMC por observación de modelos SBM y LMV para base de datos de validación 2 a) y acercamiento b). Datos Saturados superiormente (máximo 1.5).	80
Figura 4.27:	Potencia activa generada en sistema turbina-generador.	81
Figura 4.28:	EMC de las observaciones mediante modelo SBM en proceso de ajuste de variables.	83
Figura 4.29:	Variable N° 68 de la base de datos de entrenamiento.	83
Figura 4.30:	a) Vector de máximos pesos de ω y b) EMC de las observaciones según Modelo SBM Final. -	85

Figura 4.31: Resultados modelo SBM final, estadístico de Hotelling para las estimaciones.....	86
Figura 4.32: Error de estimación de las variables de salida del modelo SBM final para base de datos de entrenamiento.....	86
Figura 4.33: Error de estimación de variables de salida de modelo SBM final, base de datos T2.....	88
Figura 4.34: Indicador de anomalías según modelo SBM final, base de datos T2.....	89
Figura 4.35: Error de estimación de variables de salida de modelo SBM actualizado, base de datos T2.....	90
Figura 4.36: Indicador de anomalías según modelo SBM actualizado, base de datos T2.....	91
Figura 4.37: Comparación entre potencia generada e indicador de anomalías para modelo SBM actualizado.....	92
Figura 4.38: Acercamiento a zona de interés de Figura 4.37.....	92

ÍNDICE DE TABLAS

Tabla 3.1: Paquete de bloques para identificación de procesos industriales usando modelos SBM.....	39
Tabla 4.1: Bases de datos utilizadas.....	60
Tabla 4.2: Conjunto inicial de variables de entrada y salida.....	63
Tabla 4.3: Resumen de resultados obtenidos a lo largo de ejecución de algoritmo de generación automática de modelos SBM.....	75
Tabla 4.4: Resultados de modelos SBM y LMV para base de datos T1 según división pares/impares.....	77
Tabla 4.5: Resultados de modelos SBM y LMV para base de datos T1 según división 1-3460 y 3461-6919.....	78
Tabla 4.6: Conjunto inicial de variables de entrada y salida.....	82
Tabla 4.7: Listado de variables de entrada y salida utilizadas en modelo SBM final.....	84

CAPÍTULO 1. INTRODUCCIÓN

1.1 MOTIVACIÓN

La detección de anomalías es un problema importante que ha sido explorado en diversos campos de investigación y áreas de aplicación. Los métodos de detección de anomalías actuales permiten estudiar las irregularidades encontradas en datos históricos de un proceso, lo cual se traduce en información práctica significativa (y a veces crítica) en una amplia variedad de aplicaciones, siendo fundamental en el contexto de reducción de costos —tanto humanos como económicos— en la industria contemporánea. Es así como, es común encontrar herramientas de detección implementadas en industrias relacionadas con la minería, construcción, electricidad, transporte, o cualquier otra que involucre la manipulación de maquinaria pesada; pues una falla sostenida en estos sistemas probablemente desencadenará en un evento catastrófico que será de alto riesgo para los operadores, o bien para los propios receptores del servicio, causando pérdidas no sólo materiales sino que humanas. Por otro lado, cuando las fallas mecánicas no son identificadas —y aisladas— a tiempo, en general causan daños colaterales en la maquinaria, obligando a la industria a detener la producción para resolver el problema. En cambio, si las posibles fallas se identificaran a tiempo, sería posible prevenir el evento catastrófico, deteniendo la producción sólo para las realizar acciones de mantenimiento establecidas.

En cuanto a los mecanismos utilizados para realizar la identificación de condiciones de operación anormales o de falla, el crecimiento sostenido de la capacidad computacional ha permitido ocupar técnicas numéricas que décadas atrás sólo presentaban un enfoque teórico imposible de implementar en línea. Particularmente, los recursos tecnológicos existentes en estos días permiten la utilización de enfoques de modelación no-paramétricos los que, si bien no dependen del conocimiento fenomenológico de los sistemas dinámicos a analizar, tienen que procesar una cantidad considerable de datos cada vez que nuevas observaciones están disponibles.

1.2 ALCANCE

La presente Memoria de Título está focalizada a desarrollar una herramienta de detección de anomalías en procesos industriales multivariados, basada en la técnica de modelación no paramétrica *Similarity-Based Modeling (SBM)*. La herramienta consiste principalmente de un sistema de generación automática de modelos SBM, una metodología de estudio de eventos y un algoritmo estadístico de detección. Por lo tanto, el desarrollo del trabajo aquí presentado estará enfocado principalmente a implementar tales sistemas. Está considerado dentro del trabajo, la realización de pruebas de concepto y experimentales del sistema creado sobre un proceso industrial multivariado. Para esto, se dispone de un conjunto de bases de datos históricos de la operación de una planta de generación termoeléctrica de ciclo combinado, proporcionados por la empresa ENDESA-CHILE.

1.3 OBJETIVOS

El objetivo principal de este trabajo es la creación de una herramienta de detección de anomalías en procesos industriales multivariados usando modelos basados en similitud.

Para lograr tal propósito se han planteado una serie de objetivos específicos que permiten, en su conjunto, realizar tal tarea. Estos objetivos representan los hitos más importantes en los cuales se desarrolla el trabajo, ellos son:

- Diseño de algoritmo de generación automática de modelos SBM, basado en la selección de observaciones representativas de un proceso industrial multivariado, tanto en condición normal como anormal.
- Estudio de un proceso industrial y de sus variables significativas.
- Creación de una metodología de estudio de eventos anómalos de proceso industrial.
- Diseño de una herramienta de detección en línea de anomalías en procesos industriales usando modelos basados en similitud.

1.4 INDICACIÓN SOBRE CONFIDENCIALIDAD

El trabajo de título está inserto en un proyecto de la Empresa CONTAC INGENIEROS LTDA. en colaboración con la empresa ENDESA-CHILE. CONTAC es una empresa dedicada principalmente a servicios de Ingeniería en la aplicación de tecnologías de automatización, informática industrial y comunicaciones. Debido a la presencia de contratos de confidencialidad entre las distintas partes que constituyen el proyecto, existe información relevante referida principalmente a los algoritmos desarrollados y las bases de datos utilizadas que no será entregada de forma íntegra y detallada. Sin embargo, se deja constancia en este documento que el profesor guía de la Memoria de Título está en conocimiento de toda la información no presentada.

1.5 ESTRUCTURA GENERAL

El trabajo está constituido por 5 capítulos. En los primeros capítulos se realiza una introducción, se presenta la revisión bibliográfica y el estado del arte sobre las herramientas de detección de anomalías actuales. Se detallan los distintos enfoques posibles para abordar el problema de detección y se presenta un marco teórico sobre la detección de anomalías basada en observadores; modelación no paramétrica y análisis estadístico multivariable, conceptos utilizados en el desarrollo de este trabajo.

El Capítulo 3 describe las herramientas de modelación y detección de anomalías implementadas. En primer lugar, se describe la forma en que se utilizan técnicas estadísticas multivariadas en la creación de modelos SBM, se detalla el algoritmo de generación automática de modelos SBM y se explica la metodología de estudio de eventos elaborado. Por último, se presenta una sección que unifica los sistemas antes explicados para dar forma a la herramienta de detección de anomalías propuesta.

El Capítulo 4 presenta la descripción de un proceso industrial existente, el cual es utilizado para probar y validar la herramienta descrita en el capítulo anterior. Además se entregan los resultados obtenidos en dicho proceso, los cuales son enseñados de forma separada para cada una de las

etapas que constituyen la herramienta. Análisis de tales resultados son también presentados en este capítulo.

Finalmente, el Capítulo 5 entrega las conclusiones finales del trabajo realizado, además de proponer recomendaciones para labores futuras.

CAPÍTULO 2. REVISIÓN BIBLIOGRÁFICA

En los procesos industriales existe una gran preocupación en pos de fabricar productos de mayor calidad, de aumentar la eficiencia, reducir el número de tasas de rechazos y satisfacer normativas ambientales y de seguridad que son cada vez más exigentes. Para alcanzar tales niveles de exigencia, la gran mayoría de los procesos han tenido que aumentar el número de variables controladas. Estos controladores (PID, controladores predictivos, etc.) son diseñados para actuar bajo perturbaciones y pequeños cambios que pueda sufrir el proceso; sin embargo, existen cambios que pueden reducir la efectividad de tales controladores. Ejemplos de estas situaciones en un sistema industrial incluyen cambios en los parámetros del proceso, cambios en las perturbaciones, problemas en actuadores y problemas en sensores. Para asegurar que el proceso cumpla con las especificaciones de desempeño es necesario detectar y diagnosticar estos cambios. Estas tareas son asociadas a *monitoreo de procesos*.

El presente capítulo comienza con la Sección 2.1 la cual entrega los conceptos básicos relacionados a monitoreo de procesos y diagnóstico de anomalías. La Sección 2.2 introduce las técnicas de detección de anomalías utilizadas en la actualidad. Finalmente, la Sección 2.3 describe la detección de anomalías basado en residuos, la cual consta de la *generación de residuos* y el *análisis de residuos*. En la primera parte de la sección se presentan distintos tipos de modelos haciendo énfasis en los modelos lineales en los parámetros y no paramétricos como SBM. En la segunda parte, se describen distintas herramientas estadísticas multivariadas.

2.1 TÉCNICAS DE MONITOREO DE PROCESOS: CONCEPTOS BÁSICOS.

Para definir un vocabulario común, se presenta a continuación una serie de definiciones generadas en consenso por académicos e investigadores a nivel internacional y las cuales pueden ser encontradas en forma concisa en [24], [25] y [5].

- **Anomalía.** Son patrones en los datos medidos en un proceso que no se ajustan a un concepto bien definido de comportamiento normal. Estos patrones no conformes se denominan como anomalías, valores atípicos, observaciones

discordantes, excepciones, aberraciones, sorpresas, peculiaridades o contaminantes en los diferentes dominios de aplicación.

- **Falla (*fault*).** Desviación no permitida, con respecto a lo aceptable, usual o condición nominal, de a los menos una propiedad característica o parámetro de un sistema.
- **Evento Crítico (*failure*).** Interrupción permanente de la capacidad de un sistema para realizar una función requerida, bajo condiciones de operación específicas.
- **Mal funcionamiento.** Irregularidad intermitente en la realización de una función deseada de un sistema.
- **Falla abrupta.** Caracterizada como una función escalón. Corresponde a una falla severa que ocurre instantáneamente, por ejemplo pérdida de un sensor, bloqueo de un actuador o desconexión de una componente. Puede representar el “sesgo” en una señal medida.
- **Falla incipiente.** Caracterizada por una función rampa. De lenta evolución, por ejemplo envejecimiento o filtración. Puede representar una tendencia creciente o decreciente en una señal medida.
- **Síntoma.** Cambio en una variable observada, respecto al valor nominal.
- **Perturbación.** Una entrada que actúa sobre un sistema, la que resulta en una desviación temporal de una condición actual.
- **Residuo.** Indicador de anomalía, basado en la desviación entre valor medido y valor calculado (generalmente empleando un modelo).

El principal objetivo de las técnicas de monitoreo o supervisión de procesos es asegurar la operatividad del proceso reconociendo, en forma anticipada, anomalías en el comportamiento observado [10]. La información disponible no solo mantiene a los operadores de planta y personal de mantenimiento mejor informados con respecto al estado del proceso, sino que también, permite ayudarlos a elegir más apropiadamente las acciones que corregirán eventuales comportamientos anormales. Como resultado de este monitoreo, se minimizan los tiempos donde el proceso está detenido y sus costos económicos asociados, se mejora la seguridad en la operación de la planta, y se reducen los costos de producción.

A medida que los sistemas industriales se vuelven más integrados y complejos, la detección de anomalías que se presentan resultan ser desafíos aún más difíciles de superar si se utilizan las técnicas univariadas como *Shewhart Chart* [6], suma acumulada (CUSUM, del inglés: Cumulative Sum) [37] y media móvil con ponderación exponencial (EWMA, del inglés:

Exponentially Weighted Moving Average) [28], [16]. Estos se han diseñado para sistemas de menor escala. La debilidad de las técnicas univariadas en procesos multivariados ha obligado a realizar un gran esfuerzo en la comunidad académica e industrial en la investigación y desarrollo de técnicas de monitoreo multivariado [10]. El crecimiento en la investigación de esta índole se debe a que las industrias modernas se diseñan con una mayor cantidad de instrumentos que, naturalmente, producen una gran cantidad de información disponible. Además, el desarrollo de la tecnología ha permitido utilizar computadores cada vez más poderosos en procesamiento y superiores en capacidad de almacenamiento, permitiendo guardar un mayor cantidad de datos durante condiciones normales y anormales de los procesos.

Existen 3 procedimientos asociados al diagnóstico de una anomalía/falla: Detección de anomalía, aislamiento de anomalía e identificación de anomalía, todos enfocados a realizar una posterior recuperación (acomodación) del proceso (ver Figura 2.1) [25]. En la primera etapa se determina cuándo se ha producido una anomalía, es decir, se realiza la tarea de encontrar patrones en los datos que no se ajustan a la conducta que se espera. La segunda etapa define el tipo, localización e instante de detección de la anomalía, el propósito principal de esta fase es dirigir la atención del operario y/o ingeniero del proceso en la zona particular que merece observación. La tercera fase precisa el tamaño y (si corresponde) el comportamiento variante en el tiempo de la anomalía. En términos generales, el diagnóstico de anomalías consta de los primeros 3 puntos mencionados, es así como se determina: el tipo, ubicación, magnitud y tiempo de la anomalía. Por último, la cuarta etapa tiene que ver con la intervención del proceso con el objetivo de remover la anomalía y con este último se cierra el ciclo. El objetivo del trabajo de título se centra mayormente en la primera etapa y en menor grado a la segunda etapa de este ciclo. La literatura ofrece diversos métodos que proveen una manera de medir y tratar los datos disponibles del proceso, de tal manera de obtener información valiosa que permita guiar al operador sobre el estado del mismo. Estos mecanismos serán mencionados a lo largo de este capítulo.

De acuerdo a lo introducido con anterioridad, los esfuerzos de este trabajo están enfocados a la primera y segunda etapa del esquema de la Figura 2.1.

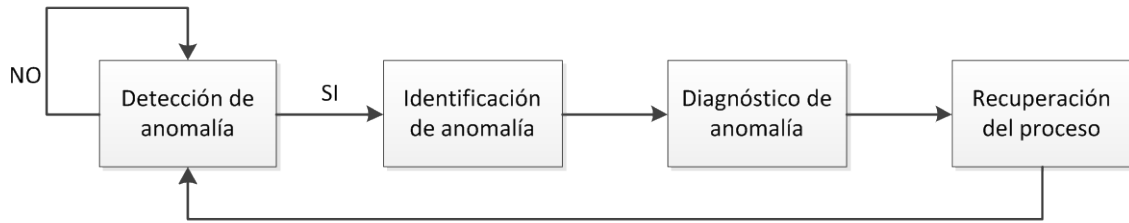


Figura 2.1: Esquema de ciclo diagnóstico de anomalías/fallas y recuperación (acomodación) del proceso.

2.2 TÉCNICAS DE DETECCIÓN DE ANOMALÍAS.

Existen diversos métodos de detección de anomalías, la mayoría enfocados a análisis en línea de la información disponible. Todos ellos se basan en definir una región que representa un comportamiento normal y declarar cualquier observación de los datos que no pertenecen a esta región normal como anormal; sin embargo, varios factores hacen de esta tarea aparentemente sencilla en un interesante desafío [5], entre estos:

- Definir una región normal, que abarca todos los comportamientos normales posibles no es una tarea fácil. Además, la frontera entre el comportamiento normal y anormal a menudo no es precisa.
- En muchos ámbitos el comportamiento normal sigue evolucionando y un concepto actual de la conducta normal podría no ser suficientemente representativo en el futuro.
- Por lo general, no se dispone de datos correctamente etiquetados y válidos para entrenamiento / validación de los modelos utilizados por las técnicas de detección de anomalías.
- Frecuentemente, los datos contienen ruido que tiende a ser similar a las anomalías reales.

Debido a los problemas mencionados, la detección de anomalías es por lo general, un problema difícil de resolver. De hecho, la mayoría de las técnicas de detección de anomalías existentes resuelven una formulación específica del problema. La formulación es inducida por diversos factores como la naturaleza de los datos y por ende del proceso; la disponibilidad de los datos y el tipo de anomalías a detectar. Es así como, por un lado existen métodos que hacen énfasis en la detección y clasificación de patrones, utilizando técnicas estadísticas multivariadas como el

análisis de componentes principales (PCA, del inglés: Principal Component Analysis) y análisis basado en el discriminante de Fisher (FDA, del inglés: Fisher's Discriminant Analysis), los cuales permiten detectar, de forma oportuna, anomalías en procesos industriales basados en patrones estadísticos obtenidos de las variables del proceso [14], [4], [18]. El uso de la transformada de Fourier o la transformada Wavelet permiten también, extraer características de las variables del proceso los cuales desembocan en técnicas exitosas de detección de anomalías [41], [21]. Otro tipo de enfoques de detección es la utilización de sistemas expertos, los cuales tienen como objetivo imitar el razonamiento humano en la detección y posterior diagnóstico de una anomalía. Los sistemas expertos están diseñados para capturar las decisiones y asociaciones que los humanos realizan, las cuales son difíciles de llevar a modelos matemáticos o causales. En la industria actual existen diversos sistemas de detección y diagnóstico de anomalías basado en sistemas expertos, desarrollados la mayor parte de ellos, por los mismos ingenieros y técnicos a cargo, los cuales han generado buenos resultados [19], [22], [15], [44].

Otra de las técnicas de detección de anomalías, ampliamente utilizada en la industria y que caracteriza el método propuesto en este trabajo, es el basado en residuos. Esta técnica puede ser descrita como la determinación de las anomalías de un sistema mediante la comparación entre mediciones disponibles del mismo y la información a priori representada por un modelo matemático. Las señales indicadoras de anomalías son denominadas *residuos*, los cuales son independientes del punto de operación del sistema y responden a las anomalías en forma característica. Este tipo de técnicas es abordado detalladamente a continuación.

2.3 DETECCIÓN DE ANOMALÍAS BASADO EN RESIDUOS

El proceso de detección de anomalías basado en residuos consta de 2 etapas: *Generación de residuos y análisis de residuos* [35]. La generación de residuos es el procedimiento en el cual se extraen síntomas de anomalía, usando información disponible de las entradas y salidas del sistema monitoreado. Los síntomas de la anomalía son representados por las señales residuales, las cuales son normalmente cercanas a cero; desviándose de este valor en forma característica cuando ocurre una anomalía. El análisis de residuos es la etapa en la cual se estudian los residuos, con el propósito de extraer información de ellos y evaluarla. Aquí se aplica una regla de decisión al residuo o una función de decisión para determinar si existe o no una anomalía, y determinar su

origen. La forma más simple de análisis es la utilización de umbrales predefinidos que, en los instantes en que se superen, permiten indicar que están ocurriendo comportamientos anormales en el proceso. A continuación se procede a entregar una descripción detallada de las dos etapas mencionadas.

2.3.1 MÉTODOS DE GENERACIÓN DE RESIDUOS

Las tres principales metodologías de generación de residuos son espacio de paridad, estimación de parámetros y basado en observador.

- Estimación de parámetros: En este método, los residuos se generan de la diferencia entre parámetros de modelos nominales y parámetros de modelos estimados. Las desviaciones en los parámetros del modelo son la base para detectar anomalías, algunos trabajos que utilizan este enfoque son [3], [26] y [23].
- Relaciones de paridad: Esta técnica consiste en formar ecuaciones a partir del modelo sin anomalías, las cuales son válidas si el funcionamiento nominal del sistema no se ve alterado por las mismas. Las relaciones de paridad son sometidas a una transformación dinámica lineal, donde los residuos transformados se usan para la detección de anomalías, como se presenta en [12], [45] y [27].
- Observadores: En esta técnica, se utiliza un observador como generador de residuos, y el residuo es generado por el error de estimación de las salidas del sistema a monitorear, es decir:

$$r(t) = f(y(t) - \hat{y}(t)) \quad (2.1)$$

Con $y(t)$ vector de variables de salida medido en planta e $\hat{y}(t)$ vector de variables de salida estimado por observador. Entonces, a diferencia del empleo del observador para fines de control, donde lo que se requiere es estimar los estados no medibles; cuando éste es utilizado para generar residuos, lo importante es estimar los estados medibles. En la detección de fallas lo importante es estimar las salidas del sistema, y no necesariamente todos los estados [10], algunos ejemplos de la aplicación de este tipo de técnicas son [17] y [11].

Este trabajo propone utilizar la metodología basada en observadores, la cual provee diversas herramientas que permiten construir modelos de un sistema.

2.3.2 GENERACIÓN DE RESIDUOS BASADO EN OBSERVADORES

El funcionamiento de un sistema de detección de anomalías basado en observadores depende en gran parte de la representación del sistema. La elección de un modelo u otro dependerá de las características del sistema a modelar y, de acuerdo a las distintas propiedades que puede presentar un modelo, se pueden distinguir las siguientes clasificaciones [2]:

- **De caja blanca, de caja negra y caja gris** [39]. Los modelos fenomenológicos o de caja blanca se construyen en base a ecuaciones físicas del sistema; los modelos empíricos o de caja negra se construyen a partir de datos de entrada y salida; los modelos grises son similares a los de caja negra, pero incorporan parcialmente aspectos fenomenológicos.
- **Continuos y discretos** [34]. Los modelos continuos y discretos son aquellos cuyas variables medidas son señales de tiempo continuo y tiempo discreto, respectivamente.
- **Modelos estáticos y dinámicos** [34]. Un modelo es estático si su comportamiento depende sólo de los valores de sus entradas en el instante actual. Por otra parte, un modelo es dinámico cuando evoluciona dependiendo de sus entradas presentes y pasadas.
- **Invariantes y variantes en el tiempo** [42]. Un modelo es invariante en el tiempo si su dinámica permanece inalterada, a diferencia de un modelo variante en el tiempo, cuya dinámica evoluciona en el tiempo.
- **Lineales y no lineales** [42]. Esta clasificación es una de las más importantes en la teoría de identificación. Surge a partir de la propiedad de linealidad.
- **Lineales y No lineales en parámetros** [42]. La propiedad de linealidad en parámetros se verifica cuando la estructura de un modelo respecto a sus parámetros cumple la propiedad de linealidad.
- **Determinísticos y estocásticos** [29]. Si un modelo incorpora procesos estocásticos como parte de sus variables y/o parámetros, se dice que el modelo es estocástico. En caso contrario, se denomina determinístico.
- **Paramétricos y no paramétricos** [33]. Aquellos modelos matemáticos que están completamente caracterizados por un conjunto finito de coeficientes (parámetros) se denominan paramétricos. En caso contrario, aquellos caracterizados por un conjunto infinito de coeficientes son llamados no paramétricos. En estadística, los modelos paramétricos hacen la hipótesis de que los datos observados son realizaciones de variables aleatorias con distribuciones de probabilidad definidas, y se hacen inferencias

sobre los parámetros de dichas distribuciones. Por el contrario, los modelos no paramétricos no realizan inferencias estadísticas sobre los datos observados, diferenciándose de los modelos paramétricos en que la estructura del modelo no se especifica a priori, sino que se determina a partir de los datos.

- **Otras características.** Otras clasificaciones existentes son: causales y no causales, modelos de parámetros concentrados o distribuidos [29].

El interés de este trabajo está enfocado al análisis según la clasificación de modelos paramétricos y no paramétricos. Como se mencionó anteriormente, esta clasificación se basa en el análisis estadístico de los datos observados.

La construcción de modelos paramétricos se basa en el supuesto de que los datos provienen de un tipo de distribución de probabilidad y hace inferencias acerca de los parámetros de la distribución. Algunos modelos hacen uso de funciones de transferencia, ecuaciones diferenciales o ecuaciones de diferencias. La metodología de construcción consta de una selección de la estructura del modelo, una formulación de criterios, estimación de parámetros y validación del modelo obtenido. Dentro de la gran cantidad de modelos que entran en esta clasificación se encuentran algunos tipos de modelos no lineales ANN (Artificial Neural Networks), los modelos lineales estáticos, lineales y no lineales de mínimos cuadrados y los modelos dinámicos estocásticos BJ, ARX, ARMAX, ARIX, ARIMAX [30].

2.3.2.1 MODELOS PARAMÉTRICOS

Los modelos paramétricos están caracterizados por un conjunto finito de coeficientes y están basados en supuestos sobre los tipos de distribución de probabilidad de los datos. A continuación, se presentan dos teorías de modelamiento paramétrico que basan su construcción en la estimación de parámetros, estos son: modelación lineal en los parámetros y regresión por mínimos cuadrados parciales (PLS).

2.3.2.1.1 MODELOS LINEALES-EN-LOS-PARÁMETROS

2.3.2.1.1.1 CASO MISO (MULTIPLE INPUTS SINGLE OUTPUT)

Los modelos entrada-salida que se utilizan para estimación de parámetros presentan en general dos entradas: la variable manipulada u y una perturbación e que representa simultáneamente perturbaciones no medidas y errores de modelación. Por simplicidad, generalmente se considera que e es un ruido blanco gaussiano de varianza desconocida.

Un modelo matemático es descrito como una relación funcional entre variables. En particular, aquellos modelos que relacionan un conjunto de variables de entrada con uno de variables de salida. Para esto es conveniente pensar en la forma genérica, es decir, la respuesta de una variable de salida Y depende de una o más entradas. X_1, X_2, \dots, X_k . Cuya relación puede ser definida como:

$$Y = f(X_1, X_2, \dots, X_k) \quad (2.2)$$

La variable Y es denominada *dependiente, respuesta o endógena* mientras que las variables X se denominan *independientes, predictores o regresores*. Esta relación permiten predecir los valores de las respuestas (a partir de los regresores); determinar el efecto de cada predictor sobre la respuesta y; confirmar, sugerir o refutar relaciones teóricas.

Los modelos lineales en los parámetros con k variables de entrada presentan la siguiente relación:

$$Y = \sum_0^k \beta_j X_j \quad (2.3)$$

Además, dada la naturaleza de los fenómenos modelados, es necesario introducir una componente aleatoria procedente de no incluir variables importantes, errores aleatorios, errores de medida y/o especificación incorrecta de la forma de la ecuación. Es así como el modelo lineal en los parámetros será del estilo:

$$Y = \sum_0^k \beta_j X_j + \epsilon \quad (2.4)$$

Donde ϵ es el error o perturbación aleatoria y los coeficientes β_j son los parámetros estructurales o estructura paramétrica de la relación propuesta.

El caso general es representado por elementos matriciales, pues se dispone de n observaciones que serán utilizadas para estimar los parámetros desconocidos:

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix} \quad (2.5)$$

Donde se incluye una columna de unos para tener en cuenta el término independiente del modelo.

El modelo es para cada una de las observaciones:

$$\begin{aligned} y_1 &= \hat{\beta}_0 + \hat{\beta}_1 x_{11} + \hat{\beta}_2 x_{12} + \cdots + \hat{\beta}_k x_{1k} + e_1 \\ y_2 &= \hat{\beta}_0 + \hat{\beta}_1 x_{21} + \hat{\beta}_2 x_{22} + \cdots + \hat{\beta}_k x_{2k} + e_2 \\ &\quad \vdots \\ y_n &= \hat{\beta}_0 + \hat{\beta}_1 x_{n1} + \hat{\beta}_2 x_{n2} + \cdots + \hat{\beta}_k x_{nk} + e_n \end{aligned} \quad (2.6)$$

Escrito de forma matricial:

$$y = X\hat{\beta} + e \quad (2.7)$$

Que corresponde al modelo poblacional:

$$y = X\beta + \epsilon \quad (2.8)$$

Con:

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_0 \\ \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix} \quad \hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} \quad e = \begin{pmatrix} e_0 \\ e_1 \\ \vdots \\ e_n \end{pmatrix} \quad (2.9)$$

El modelo anterior considera los siguientes supuestos [38]:

- **Relación estocástica:** El término de error ϵ_i recoge la componente aleatoria de y_i que el modelo no puede explicar.
- **Ausencia de error de especificación:** Aparecen en el modelo todas las variables relevantes para explicar el comportamiento de Y .
- **Linealidad de la relación:** $E(y) = X\beta$.
- **Esperanza nula del termino de perturbación:** La especificación correcta del modelo hace que no se introduzca ninguna componente sistemática en los errores al compensarse, en promedio, los positivos y negativos.
- **Homocedasticidad:** varianza constante de los errores: $Var(\epsilon_i) = \sigma^2, \forall i$.
- **No auto-correlación:** Ausencia de covarianza entre los errores: $Covar(\epsilon_i, \epsilon_j) = 0$, si $i \neq j$.
- **Variables explicativas deterministas:** Variables medidas sin error.
- **No multicolinealidad:** Las variables explicativas no son linealmente dependientes (ninguna de ellas puede obtenerse como combinación lineal de las demás).
- **Parámetros invariantes en el tiempo:** Se asume una única estructura válida para el periodo de observación y el horizonte de predicción.
- **Normalidad:** Los errores siguen una distribución normal.

2.3.2.1.1.2 CASO MIMO (MULTIPLE INPUTS MULTIPLE OUTPUTS)

El primer caso visto considera un modelo lineal de salida univariable, esto es, la respuesta y es escalar, o equivalentemente y es un vector columna si se consideran n observaciones. Los conceptos que han sido descritos para el caso antes visto pueden ser fácilmente extendidos al caso en donde la respuesta es una matriz Y , de tamaño $n \times q$, con q número de variables dependientes. En tal caso, el modelo continúa siendo el descrito en (2.3), donde:

$$Y = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1q} \\ y_{21} & y_{22} & \dots & y_{2q} \\ \vdots & \vdots & \dots & \vdots \\ y_{n1} & y_{n2} & \dots & y_{nq} \end{pmatrix} \quad (2.10)$$

Y:

$$\beta = \begin{pmatrix} \beta_{01} & \beta_{02} & \dots & \beta_{0q} \\ \beta_{11} & \beta_{12} & \dots & \beta_{1q} \\ \vdots & \vdots & \dots & \vdots \\ \beta_{k1} & \beta_{k2} & \dots & \beta_{kq} \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_{01} & \epsilon_{02} & \dots & \epsilon_{0q} \\ \epsilon_{11} & \epsilon_{12} & \dots & \epsilon_{1q} \\ \vdots & \vdots & \dots & \vdots \\ \epsilon_{n1} & \epsilon_{n2} & \dots & \epsilon_{nq} \end{pmatrix} \quad (2.11)$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_{01} & \hat{\beta}_{02} & \dots & \hat{\beta}_{0q} \\ \hat{\beta}_{11} & \hat{\beta}_{12} & \dots & \hat{\beta}_{1q} \\ \vdots & \vdots & \dots & \vdots \\ \hat{\beta}_{k1} & \hat{\beta}_{k2} & \dots & \hat{\beta}_{kq} \end{pmatrix} \quad e = \begin{pmatrix} e_{01} & e_{02} & \dots & e_{0q} \\ e_{11} & e_{12} & \dots & e_{1q} \\ \vdots & \vdots & \dots & \vdots \\ e_{n1} & e_{n2} & \dots & e_{nq} \end{pmatrix}$$

2.3.2.1.2 ESTIMACIÓN DE PARÁMETROS EN MODELOS LINEALES-EN-LOS-PARÁMETROS.

Una vez escogida la estructura del modelo, como los órdenes de cada polinomio, es necesario determinar el valor de los parámetros del mismo que ajustan la respuesta del modelo a los datos de entrada-salida experimentales. Es importante destacar, sin embargo, que esta etapa del proceso de generación se ve facilitada por la existencia de herramientas computacionales que proporcionan diferentes algoritmos para el ajuste de parámetros. Una de estas herramientas es el Tolbox de identificación de MATLAB®.

2.3.2.1.2.1 ESTIMACIÓN DE PARÁMETROS VÍA MÍNIMOS CUADRADOS (LS, LEAST SQUARE)

Existen varios métodos o criterios para realizar el ajuste de parámetros, entre los que cabe destacar el método de mínimos cuadrados y el de variables instrumentales, ambas descritas en [30].

El método de mínimos cuadrados utiliza la expresión de error de predicción:

$$e(t, \theta) = y(t) - y_e(t, \theta) \quad (2.12)$$

Con $y(t)$ salida observada del sistema, $y_e(t, \theta)$ es la salida estimada por el modelo en el instante t y θ es el vector de parámetros del modelo. Si la estructura posee regresión lineal entonces:

$$y_e(t, \theta) = \varphi^T(t) * \theta \quad (2.13)$$

Con $\varphi^T(t)$ un vector columna formado por las salidas y entradas anteriores.

Por lo tanto, juntando Ecuación (2.13) con Ecuación (2.12) se tiene:

$$e(t, \theta) = y(t) - \varphi^T(t) * \theta \quad (2.14)$$

Se presenta la siguiente función de error:

$$V_N(\theta) = \frac{1}{N} \sum_{t=1}^N \frac{1}{2} (y(t) - \varphi^T(t) * \theta)^2 \quad (2.15)$$

Conocida como *criterio de mínimos cuadrados* para una regresión lineal. Existe un valor de θ que minimiza la función anterior y que constituye la estimación del modelo por mínimos cuadrados:

$$\theta_{LSE} = sol \left\{ \frac{1}{N} \sum_{t=1}^N \varphi^T(t) * (y(t) - \varphi^T(t) * \theta) = 0 \right\} \quad (2.16)$$

Para este vector de parámetros, la función de error V_N toma su valor mínimo, siendo la función de pérdidas del modelo estimado.

En el método de variables instrumentales, el vector de parámetros debe cumplir la relación:

$$\theta_{VI} = sol \left\{ \frac{1}{N} \sum_{t=1}^N \xi(t) * (y(t) - \varphi^T(t) * \theta) = 0 \right\} \quad (2.17)$$

Donde los elementos del vector $\xi(t)$ son las llamadas variables instrumentales, que resultan de aplicar algún tipo de filtro lineal al vector de regresión lineal $\varphi^T(t)$. Este método es, de hecho, una generalización del método de mínimos cuadrados, que proporciona mejores resultados en aquellos casos en que existe algún tipo de correlación entre el ruido y la salida del sistema.

2.3.2.1.2.2 ESTIMACIÓN DE PARÁMETROS VÍA MÍNIMOS CUADRADOS PARCIALES (PLS, PARTIAL LEAST SQUARE)

La técnica de regresión por mínimos cuadrados parciales es conocida como PLS (Partial Least Squares Regression). PLS permite la modelación de sistemas proponiendo una reducción de la

dimensionalidad del problema, así como la maximización de la covarianza entre la matriz de datos X (variables independientes) y la matriz de datos Y (variables dependientes) para cada una de las componentes del espacio de proyección [9]. Específicamente, utiliza una matriz $X \in R^{n \times m}$ y una matriz $Y \in R^{n \times p}$, donde m es el número de variables predictores (número de variables medidas por cada observación), n es el número total de observaciones de los datos, y p es el número de variables observadas en Y.

Para poder aplicar esta reducción, los datos deben pasar por un tratamiento en el cual se centren y escalen todos los datos, es decir, a cada una de las variables independientes de X se le resta su media y se divide por su desviación estándar. Luego, la matriz X puede ser descompuesta en una matriz llamada *scores* $T \in R^{n \times a}$ y una matriz de carga $P \in R^{m \times a}$, donde a es la componente PLS (orden de reducción), más una matriz residuo $E \in R^{n \times m}$.

$$X = TP^T + E \quad (2.18)$$

La matriz TP^T puede ser expresada como la suma de productos de los vectores *scores* t_j (j ésima *columna de T*) y los vectores de carga p_j (j ésima *columna de P*).

$$X = \sum_{j=1}^a t_j p_j^T + E \quad (2.19)$$

Similarmente, la matriz Y se descompone en la matriz *scores* $U \in R^{n \times a}$, la matriz de carga $Q \in R^{p \times a}$, más la matriz residual $\tilde{F} \in R^{n \times p}$.

$$Y = UQ^T + \tilde{F} \quad (2.20)$$

La matriz UQ^T puede ser expresada como la suma de productos de los vectores *scores* u_j (j ésima *columna de U*) y los vectores de carga q_j (j ésima *columna de Q*).

$$Y = \sum_{j=1}^a u_j q_j^T + \tilde{F}. \quad (2.21)$$

Las matrices X e Y son representadas como la suma de una serie de matrices de rango uno. Si a es igual al $\min(m, n)$, entonces E y \tilde{F} son cero y PLS se reduce a la técnica de mínimos cuadrados ordinarios. Escogiendo a menor que $\min(m, n)$ se reduce el ruido y la colinealidad. El

objetivo de PLS es determinar los vectores de carga y *scores* que se correlacionan con Y describiendo una gran cantidad de variación de X.

La regresión PLS estima los vectores *scores* \hat{u}_j con el vector *scores* t_j como:

$$\hat{u}_j = b_j t_j \quad (2.22)$$

Donde b_j es un coeficiente de regresión. De forma matricial, la relación puede quedar expresada como:

$$\hat{U} = TB \quad (2.23)$$

Donde $B \in R^{a \times a}$ es la matriz de regresión lineal con $B_{jj} = b_j$, y \hat{U} tiene como sus columnas a \hat{u}_j . Juntando (2.23) y (2.20), teniendo en cuenta que se modifica la matriz residual, se tiene:

$$Y = TBQ^T + F \quad (2.24)$$

Donde F es la matriz de error de predicción. La matriz B es seleccionada tal que se minimice la norma de F, $\|F\|_2$. Los vectores *scores* t_j y \hat{u}_j son calculados para cada uno de los factores de PLS ($j = 1, 2, \dots, a$) tal que se maximice la covarianza entre X e Y para cada factor.

2.3.2.2 MODELOS NO PARAMÉTRICOS

Los modelos no paramétricos, a diferencia de los paramétricos, utilizan una menor cantidad de supuestos con respecto a los datos, estudiando la asociación entre las covarianzas y respuestas de los datos del proceso. Como este método no depende de un tipo de distribución de probabilidad en especial, es utilizado en procesos multivariados; sin embargo, requiere de un mayor gasto computacional para su construcción.

Los recursos tecnológicos existentes en estos días permiten la utilización de este enfoque, el que si bien no depende del conocimiento fenomenológico de los sistemas dinámicos, tiene que analizar una cantidad considerable de datos cada vez que nuevas observaciones están disponibles. Los métodos de modelamiento no paramétrico más conocidos y fuentes de investigación actual

son Kernel Regression (KR) [36], General Regression Neural Network (GRNN) [40], Radial basis Function Network (RBNF) [31] y Similarity based-model (SBM) [20], [43].

Este trabajo hace uso del enfoque basado en observadores y en particular en la creación de modelos no paramétricos basados en similitud, SBM. Este enfoque será comparado con técnicas de modelamiento más clásicos, como los modelos estáticos, lineales en los parámetros. Mencionado lo anterior, se presenta a continuación una descripción matemática que detalla la formulación de este enfoque de modelación.

2.3.2.2.1 MODELOS BASADOS EN SIMILITUD

Las técnicas de modelación no-paramétricas tienen la ventaja de que su implementación no necesita del conocimiento a priori de las estructuras a modelar, pues su funcionamiento radica en la identificación de relaciones entre los datos disponibles en vez de construir estructuras algebraicas con dicho datos. Un caso particular de dichos modelos es el modelo basado en similitud (SBM), el cual realiza una comparación entre datos medidos en línea y una base de datos representativa del sistema a modelar. SBM ha demostrado ser efectivo al ser usado en sistemas multivariados [43].

SBM considera el siguiente sistema:

$$y = f(x), x \in R^m, y \in R^p \quad (2.25)$$

Donde x e y son, respectivamente, la entrada y salida de un sistema estático, y $f(\cdot)$ es una función desconocida.

En el caso de tener acceso a mediciones de entrada y salida del sistema presentado, es posible definir las siguientes matrices de entrenamiento (de entrada y salida, respectivamente) como:

$$D_i = [x_1 x_2 \dots x_n] \in R^{m \times n} \quad (2.26)$$

$$D_o = [y_1 y_2 \dots y_n] \in R^{p \times n} \quad (2.27)$$

Donde:

$$y_i = f(x_i), \forall i = 1..n \quad (2.28)$$

Y los pares $[x_i, y_i]_{i=1..n}$ son una base representativa de los puntos de operación del sistema, la cual puede ser obtenida inteligentemente analizando los datos con técnicas estadísticas como PCA, PLS y test de Hotelling.

Luego, el método SBM asume que dada una entrada x^* , es posible estimar $y^* = f(x^*)$ mediante una combinación lineal de las columnas de D_o denotada por \hat{y}^* , es decir, el problema de estimar $y^* = f(x^*)$ se reduce a determinar un vector $w \in R^n$ tal que:

$$\hat{y}^* = D_o w \quad (2.29)$$

Una forma de encontrar este vector es mediante:

$$w = \frac{\hat{w}}{\Sigma \hat{w}} \quad (2.30)$$

$$\hat{w} = (D_i^T \Delta D_i)^{-1} \cdot (D_i^T \Delta x^*), \text{ con } \hat{w} \in R^n \quad (2.31)$$

Donde Δ es un operador de “similitud” [43].

SBM no está restringido a la consideración de un operador de similitud particular, sin embargo, éste debe cumplir con ciertas características. Para dos elementos $A, B \in R^u$, $A\Delta B \in R^+$ debe ser simétrico, y alcanzar su máximo en $A = B$, decayendo monótonamente con $\|A - B\|$. Un ejemplo de operador de similitud es el operador triangular, definido como:

$$A\Delta B = \begin{cases} d - \|A - B\| & \|A - B\| \leq d + \varepsilon \\ \varepsilon & \|A - B\| > d + \varepsilon \end{cases} \quad (2.32)$$

Donde $\varepsilon > 0$ es un número pequeño para asegurar $A\Delta B > 0$, y $d > 0$ es un umbral que depende de la dispersión de las muestras.

Si bien SBM es un método no-paramétrico que permite estimar la respuesta de sistemas estáticos, su aplicación puede ser fácilmente extendida a sistemas dinámicos discretos en el caso de contar con una secuencia de observaciones. Básicamente, para emular el comportamiento dinámico del sistema bajo estudio se deben incluir mediciones pasadas del sistema (tanto entradas como salidas) como regresores para predecir la salida, de esta forma el algoritmo SBM rescata la dependencia dinámica entre las variables medidas.

Consecuentemente, al disponer de un método de estimación no-paramétrico para sistemas dinámicos; de una base de datos representativa del sistema; y de mediciones obtenidas en forma secuencial, es posible utilizar el algoritmo SBM para implementar una rutina de monitoreo y detección de condiciones de operación anómalas. Esto se logra mediante comparación entre las observaciones y la predicción hecha por SBM, luego como SBM replica el comportamiento de la planta en base a condiciones de operaciones normales contenidas en la base representativa, si estas cantidades no son semejantes (en el sentido de cierto criterio preestablecido), se infiere que el sistema no está operando en una condición de operación conocida, sino que en una anómala (nuevamente, esto se tiene en el caso que la matriz de entrenamiento sea representativa de todos los modos de operación del sistema).

2.3.3 METODOLOGÍA DE ANÁLISIS DE RESIDUOS

Una vez que se tienen los residuos, se analizan para obtener información de ellos. El análisis tiene la finalidad de detectar la anomalía. El proceso de detección de anomalía puede ser establecido en términos de una función de decisión $J(r(t))$ y un valor umbral $T(t)$, siendo $r(t)$ el residuo.

Una anomalía puede ser detectada comparando la función de decisión $J(r(t))$ con una función de umbral $T(t)$, de acuerdo a la siguiente regla:

$$\begin{cases} J(r(t)) \leq T(t) \text{ para } f(t) = 0 \\ J(r(t)) > T(t) \text{ para } f(t) \neq 0 \end{cases} \quad (2.33)$$

Con $f(t)$ vector binario de anomalías no conocido.

Si el resultado de la prueba es positivo, es posible suponer que existe anomalía. Existen diversas formas de definir funciones de decisión y de determinar funciones de umbral, además para procesos multivariados complejos, el residuo generado puede ser complicado de analizar.

Los métodos tradicionales de análisis de residuos univariados, como la detección por discrepancia, consisten en establecer límites sobre los valores residuales obtenidos. En procesos multivariados, este tipo de métodos ignoran las correlaciones entre las variables (correlaciones espaciales) y las correlaciones entre diferentes medidas de una misma variable para distintas observaciones. La información en procesos multivariados está espacialmente correlacionada, porque usualmente hay una gran cantidad de lecturas de sensores durante todo funcionamiento del proceso y la variabilidad de las variables del proceso se limita a una dimensión más baja (por

ejemplo, debido al equilibrio de las leyes de conservación, como la materia o la energía). Además, los intervalos de muestreo son relativamente pequeños y la mayoría de los controladores de proceso son incapaces de eliminar todas las tendencias sistemáticas debido a los componentes inerciales, como en tanques, reactores, y sistemas de recirculación [26], [13].

La necesidad de extraer información oculta de una base de datos multivariadas y por lo tanto, manejar correlaciones espaciales en los datos monitoreados, ha convergido en un significativo aumento del desarrollo de técnicas estadísticas multivariadas [10]. Entre las técnicas más utilizadas se encuentran: Análisis de componentes principales (PCA, del inglés: Principal Component Analysis), Análisis discriminante de Fisher (FDA, del inglés: Fisher's Discriminant Analysis) y regresión de mínimos cuadrados parciales (PLS, del inglés, Partial Least Square regression). Técnicas que tienen un campo de implementación enorme y pueden ser aplicadas también a análisis de residuos para procesos multivariados.

La técnica más ampliamente utilizada en procesos multivariados es PCA. Este método reduce de forma óptima las dimensiones del proceso desde el punto de vista de la variabilidad de los datos. La estructura formada por PCA permite identificar las variables responsables de una anomalía o las variables que se ven más afectadas por la anomalía. En casos donde la mayor parte de la información importante en los datos puede ser capturada en solo 2 ó 3 dimensiones, lo cual es cierto para algunos procesos, la variabilidad dominante del proceso puede ser ilustrada en un sólo gráfico.

Independiente del número de dimensiones utilizadas, existen otros tipos de herramientas, como por ejemplo, Test de Hotelling (T_α^2), el cuál provee de un umbral escalar estadísticamente significativo para evaluar los datos multivariados. Este test es llamado estadístico de Hotelling y da cuenta de la variabilidad de los datos considerando como supuesto básico, que aquellos datos siguen una distribución normal multivariada. Esta herramienta ayuda a los operadores e ingenieros a interpretar las tendencias significativas de los datos del proceso y es indicador estadístico interesante cuando se emplea en el análisis de residuos. En donde es posible determinar la calidad de la estimación realizada por un modelo y de esa manera, utilizarlo en un algoritmo de construcción del modelo dirigiendo los esfuerzos por mejorar la estimación y optimizar la detección de posibles anomalías.

FDA es una técnica de reducción de dimensionalidad desarrollada y estudiada como metodología para clasificación de patrones. Ésta determina la porción del espacio de observación que es más eficaz en la discriminación entre clases de datos, por ello, FDA se emplea mayoritariamente como

técnica de diagnóstico de anomalías. Una de las cualidades de ésta es que se aplica a los datos de todas las clases simultáneamente. Es usada en análisis de residuos ya que permite clasificar los tipos de residuos obtenidos.

PLS es otra técnica bastante usada en sistemas multivariados ya que, como se mencionó con anterioridad, permite la modelación de sistemas proponiendo una reducción de la dimensionalidad del problema, así como la maximización de la covarianza entre la matriz de datos X (variables independientes) y la matriz de datos Y (variables dependientes) para cada una de las componentes del espacio de proyección. Debido a que la detección de anomalías depende en gran medida de la correcta identificación del proceso, es posible hacer estudios previos, basados en la técnica PLS, para determinar las variables significativas que aporten realmente a la modelación del proceso, indistintamente del tipo de modelo y del sistema/evento en estudio. Usando la reducción de dimensionalidad que provee PLS es posible disminuir significativamente el número de variables de entrada/salida de un modelo y optimizar el costo computacional sin perder calidad en la estimación de las salidas.

El presente trabajo utiliza varias de las técnicas descritas con anterioridad. La construcción de modelos SBM efectivos y robustos no es trivial y requieren de un análisis estadístico de los datos existentes, por ello, se propone un algoritmo de análisis de datos que permita construir automáticamente este tipo de observadores basado en el análisis de componente principales (PCA). Para lograr este objetivo se deberá analizar además, el residuo. El test de Hotelling será la herramienta encargada de transformar los residuos, provenientes de la diferencia entre las salidas observadas y las salidas estimadas, a residuos transformados conformándose entonces, en indicadores más confiables y robustos. En este contexto, las técnicas escogidas precisan de una explicación más detallada. Es por ello que, para entregar al lector una base teórica completa y clara de cada una de ellas, se describen a continuación.

2.3.3.1 ANÁLISIS DE COMPONENTES PRINCIPALES.

Esta técnica de reducción de dimensionalidad permite preservar la estructura de correlación entre las variables del proceso capturando de forma óptima la variabilidad de los datos [8]. Este método determina un conjunto de vectores ortogonales, llamados componentes principales, que son ordenados de acuerdo a la cantidad de varianza explicada por la dirección de cada vector. Dado un conjunto de datos con n observaciones y m variables ordenados en una matriz X .

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{pmatrix} \quad (2.34)$$

Las componentes principales se calculan al resolver el problema de optimización:

$$\max_{v \neq 0} \frac{v^T X^T X v}{v^T v} \quad (2.35)$$

Donde $v \in R^m$. El cual puede ser calculado usando la descomposición en valores singulares:

$$\frac{1}{\sqrt{(n-1)}} X = U \Sigma V^T \quad (2.36)$$

Donde $U \in R^{n \times n}$ y $V \in R^{m \times m}$ son matrices unitarias, y $\Sigma \in R^{n \times m}$ es una matriz diagonal con valores singulares reales no negativos ordenados decrecientemente ($\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{\min(m,n)} \geq 0$). Las componentes principales son los vectores columnas ortonormales de la matriz V , y la varianza del conjunto de datos proyectado en la $i^{\text{ésimo}}$ columna de V es igual a σ_i^2 . Resolver la ecuación anterior es equivalente a resolver la descomposición en valores y vectores propios de la matriz de covarianza S .

$$S = \frac{1}{n-1} X^T X = V \Lambda V^T \quad (2.37)$$

Donde la matriz diagonal $\Lambda = \Sigma^T \Sigma \in R^{m \times m}$ contiene valores reales no negativos ordenados decrecientemente de los valores propios ($\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_m \geq 0$) y el $i^{\text{ésimo}}$ valor propio equivale al cuadrado del $i^{\text{ésimo}}$ valor singular (es decir, $\lambda_i = \sigma_i^2$).

Con el fin de capturar de manera óptima la variabilidad de los datos y reducir al mínimo el efecto del ruido aleatorio que afecta a la representación PCA, se escogen las primeras a componentes principales que corresponden a los a primeros valores singulares (ordenados decrecientemente). De esta manera, se tendrá la matriz $P \in R^{m \times a}$ que contiene las primeras a componentes principales. Las proyecciones de las observaciones en X en el espacio reducido están contenidos en la matriz $T \in R^{n \times a}$ definida como:

$$T = XP \tag{2.38}$$

Y la proyección de T en el espacio de m dimensiones será:

$$\hat{X} = TP^T \tag{2.39}$$

Con esto, es posible construir la matriz residual:

$$E = X - \hat{X} \tag{2.40}$$

La matriz E captura las variaciones en las observaciones asociadas a las componentes principales que no fueron consideradas en el espacio reducido ($(m - a)$ últimas componentes). El subespacio contenido en la matriz E tiene una razón señal-ruido menor y, al ser removida del espacio X , produce una representación más precisa del proceso.

Definiendo t_i como la $i^{\text{ésima}}$ columna de T . Es posible demostrar las siguientes propiedades:

1. $Var(t_1) \geq Var(t_2) \geq \dots \geq Var(t_a)$.
2. $Promedio(t_i) = 0; \forall i$.
3. $t_i^T t_k = 0; \forall i \neq k$.
4. No existe ningún otro tipo de expansión ortogonal de a componentes que capture más variabilidad en los datos.

Un nuevo vector de observación ($x \in R^m$) es fácilmente proyectado en el espacio reducido, $t_i = x^T p_i$, donde p_i es la $i^{\text{ésima}}$ componente principal. La variable transformada t_i es llamada la $i^{\text{ésima}}$ componente principal de x . El esquema puede ser ilustrado como:

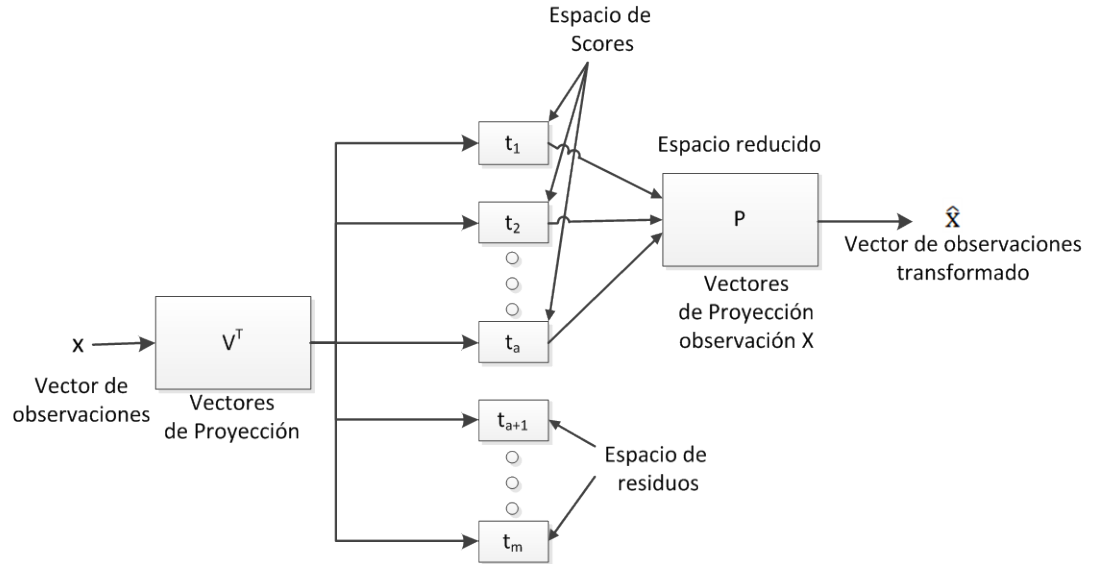


Figura 2.2: Proyección del vector de observación x , en los espacios score y de residuos, y el cálculo del vector de observación transformado.

La Figura 2.2 anterior ilustra la transformación de un nuevo vector de observación al espacio transformado, esta transformación permite representar el vector original a través de una combinación lineal de vectores de proyección definidos por PCA. La representación de tal vector es realizada con una cantidad menor de vectores de proyección, los cuales capturan la mayor variabilidad de los datos, de esta forma, se transforma el vector en su espacio original a uno en el espacio transformado reducido, el cual corresponde al vector de observación transformado por PCA.

2.3.3.2 TEST ESTADÍSTICO DE HOTELLING

Este test estadístico permite caracterizar la variabilidad de un conjunto de observaciones de variables a través de un umbral escalar, asociado a un determinado nivel de confianza. La utilidad de este tipo de test en detección de anomalías basado en modelos radica en la posibilidad de evaluar la estimación, bajo un nivel de confianza, verificando si el error de estimación se encuentra en una región aceptable para cada uno de los subconjuntos de variables [7].

Considerándose la matriz X definida en (2.34), la matriz de covarianza de X es:

$$S = \frac{1}{n-1} X^T X \quad (2.41)$$

Y su descomposición en valores y vectores propios como:

$$S = V\Lambda V^T \quad (2.42)$$

La cual revela la estructura de correlación de la matriz de covarianza, donde Λ es diagonal y V es ortogonal ($V^T V = I$). La proyección $y = V^T x$ de un vector de observación $x \in R^m$ desacopla el espacio de observación en un conjunto de variables no correlacionadas correspondientes a los elementos de y . La varianza del $i^{\text{ésimo}}$ elemento de y es igual al $i^{\text{ésimo}}$ valor propio de la matriz Λ . Asumiendo que S es invertible y la siguiente definición:

$$z = \Lambda^{-\frac{1}{2}} V^T x \quad (2.43)$$

El estadístico de Hotelling T^2 se define como:

$$T^2 = z^T z \quad (2.44)$$

La matriz V rota los ejes principales de la matriz de covarianza X de tal manera que corresponden directamente a los elementos de y , además Λ escala los elementos de y produciendo un conjunto de variables con varianza unitaria que corresponden a los elementos de z . Como ejemplo, se presenta la transformación de la matriz de covarianza para un espacio de 2 dimensiones en la Figura 2.3.

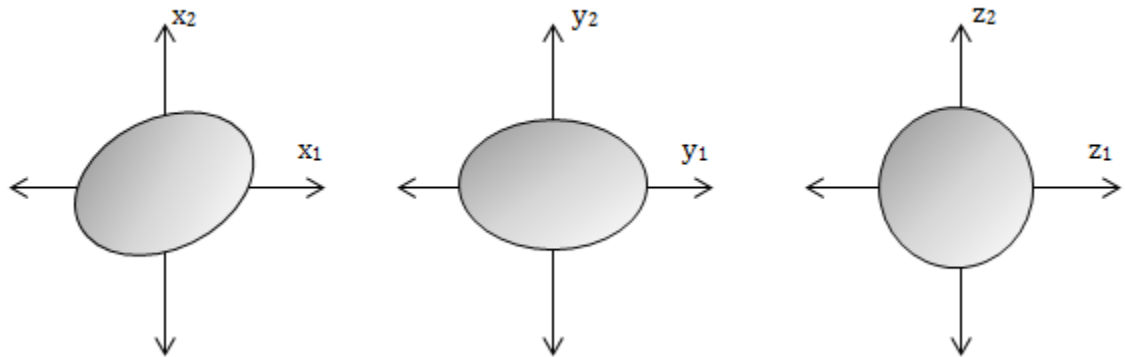


Figura 2.3: Ilustración grafica de la transformación del estadístico de Hotelling.

El umbral apropiado para el estadístico de Hotelling basado en un nivel de confianza α , puede ser determinado asumiendo que las observaciones son muestras aleatorias de una distribución normal multivariable. Considerando que el estadístico de Hotelling puede ser utilizado en dos situaciones: detectar valores atípicos en el conjunto de entrenamiento y detectar observaciones anómalas fuera de conjunto de entrenamiento. Es posible utilizarlo por ejemplo, en el análisis de la base representativa que crea modelos SBM (ver Sección 2.3.2.2.1). En esta situación, se define el siguiente umbral de Hotelling.

$$T_{\alpha}^2 = \frac{(n-1)^2 \left(\frac{m}{n-m-1} \right) F_{\alpha}(m, n-m)}{n \left(1 + \left(\frac{m}{n-m-1} \right) F_{\alpha}(m, n-m) \right)} \quad (2.45)$$

Donde $F_{\alpha}(g, k)$ es la distribución Fisher con un nivel de confianza α , con g y k grados de libertad. Además, y continuando el ejemplo, es posible analizar observaciones anómalas que se encuentran fuera del conjunto de entrenamiento del modelo SBM y así detectar anomalías en las observaciones del proceso. Para esto, se define el siguiente umbral de Hotelling:

$$T_{\alpha}^2 = \frac{m(n-1)(n+1)}{n(n-m)} F_{\alpha}(m, n-m) \quad (2.46)$$

Luego de entregado al lector una base teórica sobre los principales conceptos y métodos utilizados actualmente en el área de detección de anomalías y análisis estadístico multivariable, se presenta a continuación un capítulo que entrega una descripción del sistema de detección de anomalías usando modelos SBM propuesto, así como también, la metodología diseñada para el estudio de nuevos eventos anómalos en un proceso industrial multivariado.

CAPÍTULO 3. IMPLEMENTACIÓN DE HERRAMIENTA DE DETECCIÓN DE ANOMALÍAS USANDO MODELOS SBM.

La detección de anomalías está referida a encontrar patrones en datos que no corresponden al comportamiento esperado. Existe una gran variedad de rubros que hacen uso de sistemas de detección de anomalías, algunos de ellos son empresas financieras, actividades militares, seguridad cibernética; sin embargo, el interés de este trabajo está enfocado a los procesos industriales. La importancia de detectar anomalías en dichos procesos radica en las acciones que se deben tomar cuando ocurren tales anomalías, como son la detención total de la producción en los casos más extremos, creación de una nueva planificación de la producción y de los tiempos de mantención/renovación de equipos, o también en la forma en que debe operar un proceso, siendo necesario hacer modificaciones de tal manera de mantener un nivel de calidad/eficiencia mínimo.

El propósito de este capítulo es entregar una descripción detallada de la implementación de una herramienta de detección de anomalías basado en residuos para procesos industriales usando modelos SBM. Comienza en la Sección 3.1 proponiendo la utilización de técnicas estadísticas multivariadas en la creación de modelos SBM. La Sección 3.2 detalla el algoritmo de generación automática de modelos SBM, dando dos enfoques de distinto nivel de profundidad en la descripción. La Sección 3.3 presenta un esquema metodológico para el estudio de eventos usando la herramienta SBM generada. Finalmente, en la Sección 3.4 se entrega el procedimiento a seguir para implementar una herramienta de detección en línea de anomalías en procesos industriales multivariados.

3.1 INTRODUCCIÓN A LAS TÉCNICAS UTILIZADAS

Los modelos basados en similitud son parte de las técnicas de modelamiento no paramétrico que requieren, además de una selección apropiada de variables, tanto de entrada como de salida (abordado más adelante), de una base de datos de entrenamiento para su construcción; más aún, el

modelo SBM se genera en base a la selección de observaciones¹ puntuales que debiesen representar fielmente el proceso en sus distintos modos de operación (ver Sección 2.3.2.2.1). Esta selección es fundamental para lograr un buen desempeño del modelo.

En la actualidad, la literatura no provee de técnicas enfocadas a seleccionar de manera inteligente y automatizada tales observaciones representativas, y en general, se hace una selección manual de los instantes que, según los operadores y/o especialistas, aparentan ser los más representativos de los diferentes puntos de operación del proceso. Debido a esta razón, se propone generar un algoritmo que permitirá escoger automáticamente las observaciones más representativas de la base de datos de entrenamiento. El algoritmo, descrito más adelante, consta de varias etapas elementales las cuales hacen uso de algunas técnicas vistas en el capítulo anterior.

El mecanismo automatizado de extracción de observaciones representativas que se presenta en este trabajo está basado en 2 técnicas estadísticas multivariantes. Ellas son: análisis de componentes principales (PCA) y test estadístico de Hotelling (T_{α}^2). PCA reduce la dimensionalidad del problema, permitiendo preservar la estructura de correlación entre las variables del proceso y capturando de forma óptima la variabilidad de los datos. Esta característica hace de este método una excelente herramienta para seleccionar observaciones representativas, pues la transformación lineal que realiza PCA genera un nuevo sistema de coordenadas para el conjunto original de datos en el cual la varianza de mayor tamaño del conjunto de datos es capturada en el primer eje de proyección (Primera Componente Principal), la segunda varianza más grande es el segundo eje de proyección, y así sucesivamente (ver Figura 3.1, izquierda). De este modo es posible observar en las nuevas componentes, agrupamientos de observaciones que tienen que ver con un comportamiento similar del proceso, lo que se traduce en la determinación de observaciones que pertenecen a una misma condición de operación; es más, el número de agrupaciones de datos representa usualmente, al número de condiciones de operación disímiles. Por lo tanto, como primera selección de observaciones se propone seleccionar, por medio de un algoritmo, las observaciones que se encuentran cercanas al centro de cada una de las agrupaciones mencionadas (ver Figura 3.1, derecha) y, de esta manera, generar un modelo SBM inicial con el 4 % de las observaciones de la base de datos. Escoger solamente un 4% de las observaciones pretende evitar la inclusión de un número elevado de

¹ Se entiende por observación del proceso a la medición del valor de cada una de las variables involucradas en un proceso en un tiempo determinado.

muestras redundantes, además de reducir o simplificar los requerimientos de cómputo en operación en línea.

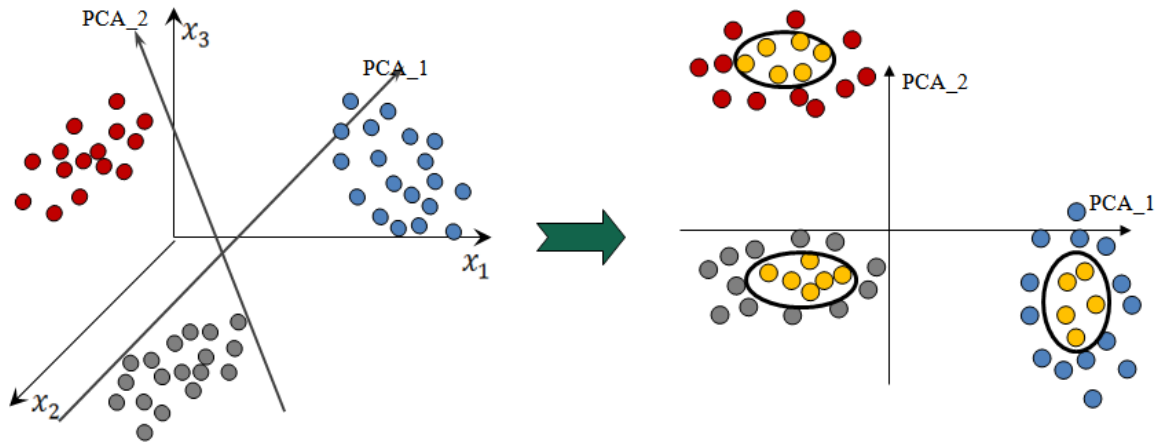


Figura 3.1: Representación de observaciones transformadas por PCA (izquierda), selección de observaciones representativas (derecha).

Como segunda utilidad de la herramienta PCA, en la que se emplea en conjunto con el estadístico de Hotelling, es en la reducción de dimensión de la matriz de errores de estimación (diferencia entre variable medida y variable estimada), con el fin de analizar un número más pequeño de variables. De esta manera, en vez de utilizar toda la matriz de errores, es posible reducir la cantidad de información, capturando de forma óptima la variabilidad de estos (reducción que captura el 95 % de la variabilidad de los datos). Esta metodología permite determinar mediante el estadístico de Hotelling con un nivel de confianza establecido, cuáles de las observaciones merecen ser consideradas en la creación del modelo dejando fuera de esta selección a las observaciones que sobrepasan el umbral de Hotelling establecido en Ecuación (2.45), pues son considerados como valores atípicos no representativos del proceso. Luego, en base a los errores obtenidos con la estimación de las variables de salida del modelo SBM inicial (4% de base de datos de entrenamiento), es posible seleccionar un 6% de las observaciones no escogidas inicialmente, para así completar una base representativa que corresponde al 10 % de los datos de entrenamiento. Esta metodología de análisis de la matriz de errores de estimación es utilizada en varias ocasiones más, donde por medio de iteraciones, se optimiza el modelo SBM usando como criterio de comparación entre modelos, el error cuadrático medio normalizado.

Con respecto a la rutina de detección de anomalías propuesto, el análisis de las estimaciones estará basado-también- en el estadístico de Hotelling. Esta vez, dicho método será encargado de determinar observaciones anómalas que se encuentran fuera del conjunto de entrenamiento del

modelo SBM y así detectar anomalías en las mediciones del proceso, basado en el umbral presentado en Ecuación (2.46).

Adicionalmente a lo ya descrito, es necesario mencionar que la metodología para el estudio de eventos que será presentado en este trabajo contiene, como parte fundamental, el análisis de variables de entrada/salida que serán utilizadas en la creación de modelos SBM. Dentro de este estudio, se propone utilizar la técnica de regresión de mínimos cuadrados parciales. Como se mencionó en el capítulo anterior, PLS es una técnica de modelación cuyo principal objetivo es explicar una o más variables dependientes (Y) en función de un número de variables explicativas (predictores, X). La idea atrás de PLS es que un gran número de variables incluidas en la matriz X captan de alguna forma, los efectos dominantes producidos por cambios en la estructura de las observaciones. Sin embargo, se debe estar consciente de que puede haber variables en la matriz X que no tienen relación con tales cambios e introducen solamente ruido en la descripción. La idea de seleccionar/eliminar variables, tanto de entrada como de salida, se basa en encontrar cuales de ellas son las que introducen mayoritariamente ruido y no están aportando a la descripción de la estructura de los datos. Por esto, se utilizará PLS para detectar variables que aumentan la capacidad predictiva de los modelos SBM. De esta manera, el procedimiento para la eliminación consta en primer lugar, de escoger un conjunto inicial de variables de entrada y salida en base a estudio previo, luego analizar los datos de tales variables con PLS y graficar los vectores de carga (ver Figura 3.2), verificando la ubicación de cada una de las variables dentro de tal gráfico, para así calcular la distancia al origen de cada uno de ellos, entendiéndose que las variables más cercanas a tal punto son exactamente las variables que no son explicadas por los vectores de carga y por lo tanto, no aportan a la correcta estimación del modelo. De esta manera, es posible eliminar variables, tanto de entrada como de salida, escogiendo un umbral apropiado en donde todas las variables cuya distancia sea menor al umbral, quedan eliminadas del modelo SBM. Esta metodología puede ser aplicada de forma iterativa, de este modo, es posible eliminar progresivamente variables al incrementar el umbral y así, evaluar distintos esquemas de modelos SBM. De esta forma y utilizando como criterio de comparación el error medio cuadrático normalizado, es posible encontrar un número de variables de entrada y salida que producen una estimación óptima en términos de la elección de variables.

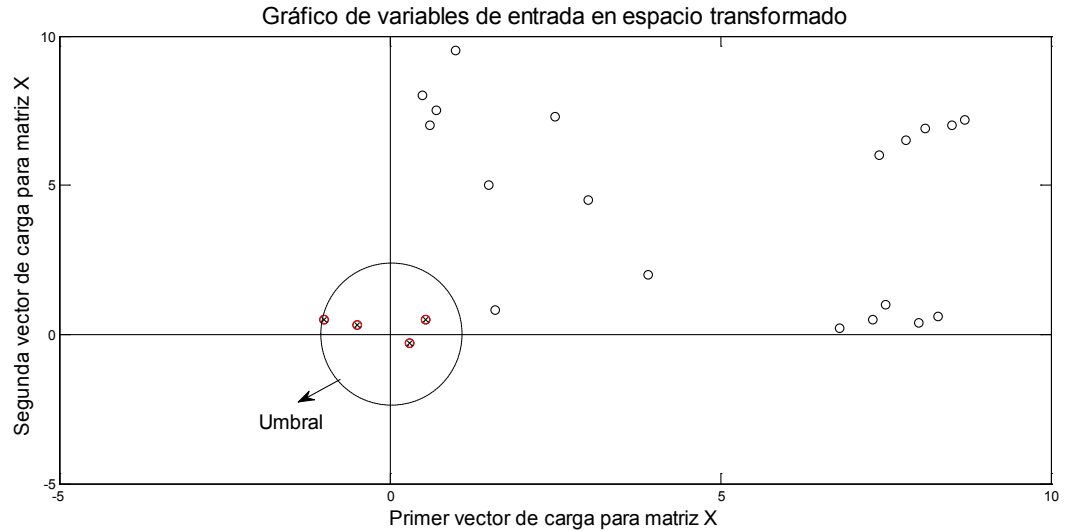


Figura 3.2: Ilustración de método de eliminación de variables basado en análisis de componentes PLS.
Representación de variables en espacio transformado.

3.2 ALGORITMO DE GENERACIÓN AUTOMÁTICA DE MODELOS SBM

Basado en la necesidad de contar con una herramienta de creación automática de modelos, que pueda ser utilizada en detección de anomalías basado en residuos, se presenta a continuación un algoritmo que construye modelos SBM de manera automatizada. El objetivo principal es que, a partir de una elección inteligente de las variables de entrada y salida a considerar en el modelo, además de una base de datos válida y rica en información de puntos de operación, se puedan generar modelos basados en similitud que permitan estimar variables de salida y así, usarlas en la herramienta de detección de anomalías propuesta en este trabajo. El algoritmo de generación automática de modelos SBM tiene como esencia principal la selección de las observaciones que serán consideradas como más representativas del proceso dentro de la base de datos de entrenamiento, siendo éstas las utilizadas para crear relaciones internas entre entradas y salidas que finalmente, permiten generar estimaciones de las variables de salida, dado un set de valores de entrada.

A modo de entregar al lector una descripción exhaustiva y completa de un algoritmo de mediana complejidad, tal como el implementado, se procederá a entregar 2 enfoques. En primer lugar se presentará una descripción superficial de lo que, a grandes rasgos, realiza el algoritmo,

representando las acciones según bloques fundamentales generales. En segundo lugar, se entrega una descripción detallada de cada uno de los pasos que se realizan para cumplir con el objetivo propuesto para el algoritmo, o sea, determinar un modelo SBM de forma automática de un proceso industrial a partir de una base de datos representativa del mismo.

El algoritmo diseñado y que se presenta a continuación fue desarrollado en su totalidad utilizando el software MATLAB®.

3.2.1 PRIMER ENFOQUE: DESCRIPCIÓN GENERAL:

Es posible fragmentar el algoritmo de creación automática de modelos SBM en los siguientes pasos (ver Figura 3.3).

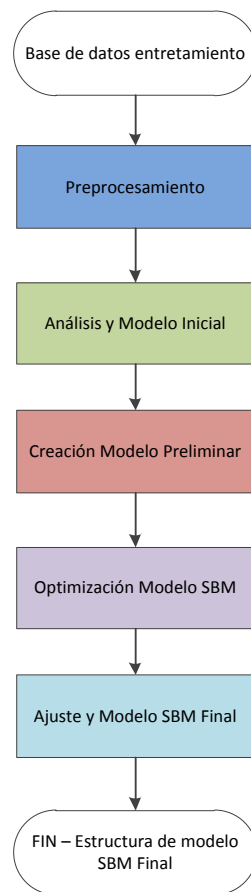


Figura 3.3: Diagrama de flujo de algoritmo - Descripción general.

Como se puede observar en la Figura 3.3, el algoritmo se inicia con una base de datos de entrenamiento. Esta base de datos debe ser representativa del proceso en estudio, rica en

información sobre condiciones de operación y lo suficientemente válida para su análisis, además se debe escoger previamente cuáles serán las variables de entrada y salida para el modelo (basados en la metodología de estudio de eventos, ver Sección 3.3). El algoritmo por lo tanto, está diseñado para leer la información desde un archivo (la forma de incorporar datos puede variar para la herramienta de detección en línea de anomalías). A modo de asegurar la robustez del algoritmo, es necesario pre-procesar los datos para eliminar posibles elementos no válidos (como por ejemplo, NaN, del inglés, Not a Number). De esta manera, se entrega a las etapas siguientes, dos matrices: $V_{inTrain} \in R^{n \times p}$ y $V_{outTrain} \in R^{n \times q}$. Donde sus columnas representan variables medidas de entrada y de salida respectivamente, y sus filas representan los instantes de tiempo en donde fueron tomados los datos (n instantes de tiempo).

Luego del pre-procesamiento se realiza un análisis inicial, el cual genera un modelo SBM a partir del 4% de las observaciones de la base de datos ingresada, utilizando para ello la técnica de análisis de componentes principales (ver Sección 2.3.3.1). Luego, al graficar las observaciones en los nuevos ejes transformados, es posible seleccionar aquellas que se encuentran en los centros de cada una de las agrupaciones formadas.

Seleccionados finalmente el 4 % de las observaciones de la base de datos de entrenamiento, es posible realizar un modelo inicial, que servirá para obtener una primera estimación de las variables de salida a través de un modelo SBM.

La estimación de las salidas obtenidas con el modelo inicial sirve para realizar un tratamiento mediante PCA y test de Hotelling a la matriz de errores, la cual es obtenida como la resta entre el valor medido en planta y el valor estimado para cada una de las variables de salida e instantes de tiempo. De esta manera y a través de un criterio de preselección basado en el estadístico de Hotelling (ver Sección 2.3.3.2), es posible incluir observaciones no consideradas en el conjunto inicial (4% de la base de datos) y así, completar un conjunto preliminar constituido por un 10 % de las observaciones de la base de entrenamiento, generando un modelo SBM preliminar que permite estimar las variables de salida y obtener una matriz de error actualizada que se utiliza en la etapa de optimización posterior.

La matriz de error generada por las estimaciones del modelo preliminar son analizados nuevamente a través de PCA y el estadístico de Hotelling, utilizando la misma función de selección de observaciones candidatas que no han sido consideradas en el modelo preliminar, es posible renovar el conjunto preliminar de observaciones y así determinar la calidad de las estimaciones generadas comparando el modelo SBM preliminar con el modelo SBM actualizado

usando como criterio el error medio cuadrático. Se propone entonces, generar 4 sub-etapas que consideran la eliminación parcial del conjunto inicial de observaciones (4%). Esta eliminación se basa en que existen observaciones dentro de tal conjunto que pueden ser redundantes (entregan misma información al modelo SBM) y/o no son representativas de las verdaderas condiciones de operación del proceso. Es por ello que, se modifica el conjunto inicial de observaciones que definen el modelo preliminar, sustituyendo un 1% de los instantes de tiempo seleccionados inicialmente por observaciones candidatas obtenidas por el último modelo SBM creado, evaluando la calidad del nuevo modelo y conservando el 1% de observaciones candidatas si es que el modelo mejora en la estimación de las salidas, en caso contrario, se reincorpora el 1% de las observaciones que fueron reemplazadas. Siguiendo este procedimiento, se realizan 4 iteraciones que permiten así ajustar el conjunto de observaciones que definirán el modelo SBM final.

Existe una etapa opcional que ajusta el modelo SBM final incorporando observaciones no consideradas por el algoritmo descrito anteriormente. El criterio de inclusión se basa en el error medio cuadrático y el estadístico de Hotelling.

Finalmente, es posible entregar una versión final del modelo SBM calculado y optimizado representado por un objeto de MATLAB® llamado *estructura* que considera los siguientes elementos:

- $D_i \in \mathbf{R}^{m \times nVin}$: Matriz con $nVin$ variables de entrada y m observaciones, siendo $m \cong 0.1 \cdot n$.
- $D_o \in \mathbf{R}^{m \times nVout}$: Matriz con $nVout$ variables de salida y m observaciones, siendo $m \cong 0.1 \cdot n$.
- $ETrain \in \mathbf{R}^{n \times nVout}$: Matriz que representa el error de estimación del modelo SBM final, para cada variable de salida y para cada una de las observaciones.
- $MeanETrain \in \mathbf{R}^{1 \times nVout}$: Vector que contiene en cada elemento el promedio aritmético de cada una de las columnas de $ETrain$.
- $SigmaETrain \in \mathbf{R}^{1 \times nVout}$: Vector que contiene en cada elemento la desviación estándar de cada una de las columnas de $ETrain$. En el caso en que una columna tenga desviación estándar igual a cero, se reemplaza tal valor por uno.

- $vpETrain \in R^{nVout \times nVout}$: Matriz que contiene los vectores propios de la matriz de covarianza de $ETrain$. Representan cada uno de las componentes principales obtenidas al hacer PCA a la matriz $ETrain$.
- $vLETrain \in R^{nVout \times 1}$: Matriz que contiene los valores propios de la matriz de covarianza de $ETrain$.

3.2.2 SEGUNDO ENFOQUE: DESCRIPCIÓN DETALLADA DE FUNCIONES:

Luego de entregado el enfoque general, se procede a explicar detalladamente el algoritmo, el cual fue fragmentado en diferentes bloques y reunidos en un paquete de identificación de procesos multivariados y ejecutables en el software MATLAB®.

3.2.2.1 BLOQUES FUNCIONALES UTILIZADOS EN IDENTIFICACIÓN DE PROCESOS INDUSTRIALES USANDO MODELOS SBM.

Los bloques implementados en MATLAB® y que permiten lograr el objetivo propuesto se resumen a continuación:

PAQUETE DE BLOQUES CREADOS EN MATLAB®

A. Pre-procesamiento
B. Limpieza de Datos
C. Análisis Inicial
D. División conjunto inicial
E. SBM
F. Similitud
G. Vectores LD
H. Error de estimación
I. Filtro mediana móvil.
J. Análisis mediante Hotelling
K. Selección Vía Hotelling
L. Ajuste Final

Tabla 3.1: Paquete de bloques para identificación de procesos industriales usando modelos SBM.

Luego de presentada la lista de bloques utilizados, se entrega a continuación una descripción de cada uno de ellos para así, explicar el algoritmo principal haciendo referencia a la letra de cada una de los bloques agrupados en la Tabla 3.1.

3.2.2.1.1 PRE-PROCESAMIENTO (A)

El bloque “Pre-procesamiento” tiene como objetivo transformar la información contenida en un archivo de datos de entrada a objetos que pueden ser manipulados en MATLAB®, como son las matrices. De esta manera, recibe como entrada: el nombre del archivo CSV donde se encuentran los datos históricos del proceso; el nombre del archivo CSV que contiene una descripción de las variables; un índice de línea de lectura inicial y final; y un vector con las variables que serán leídas.

Con esto, se leen los datos y se procesan para obtener una matriz que contiene la información del proceso; las etiquetas de las variables leídas, las fechas de las observaciones leídas y la descripción de cada variable. En el caso en que existiese información mal recuperada y/o que contenga símbolos que representen mala lectura, como por ejemplo: NaN (not a number), los

cuales estarán presentes en la matriz de salida de datos, éstos deberán ser identificados y eliminados.

3.2.2.1.2 LIMPIEZA DE DATOS (B)

El bloque “Limpieza de datos” tiene como objetivo limpiar las matrices de datos que contienen las observaciones de las variables de entrada. Principalmente, se busca eliminar elementos que no sean números. Para ello, el algoritmo realiza un barrido de cada una de las variables (por columnas), tanto de entrada como de salida, y elimina aquellas variables que superan un número determinado de elementos no válidos. Luego, se procede a realizar un barrido por observación (por filas) para eliminar cualquier observación que contenga al menos mal capturado. De esta manera, se tienen como salidas los mismos objetos de entrada sin elementos incorrectamente leídos.

3.2.2.1.3 ANÁLISIS INICIAL (C)

El bloque “Análisis inicial” tiene como objetivo generar un conjunto inicial (C_1) de observaciones representativas de la base de datos de entrenamiento. Para ello realiza un análisis de componentes principales a la base de datos de entrenamiento. Esto permite identificar puntos de operación y así, seleccionar de preliminarmente observaciones que son la base para la generación del modelo SBM inicial.

El bloque genera un conjunto C_1 que constituye el 4% de las observaciones de la base de datos de entrenamiento a partir de un análisis de componentes principales; Esta tarea es realizada transformando, mediante PCA, los datos de entrenamiento a un nuevo espacio que rescata en sus primeras 2 componentes principales la mayor variabilidad de tales datos. Con esto, es posible graficar las observaciones en los nuevos ejes transformados por PCA (2 primeras componentes principales) y seleccionar aquellas observaciones que se encuentran en los centros de cada una de las agrupaciones formadas, como se puede observar en la Figura 3.4. La determinación del número de agrupaciones y los centros de cada uno de ellos es realizado numéricamente y de forma automática.

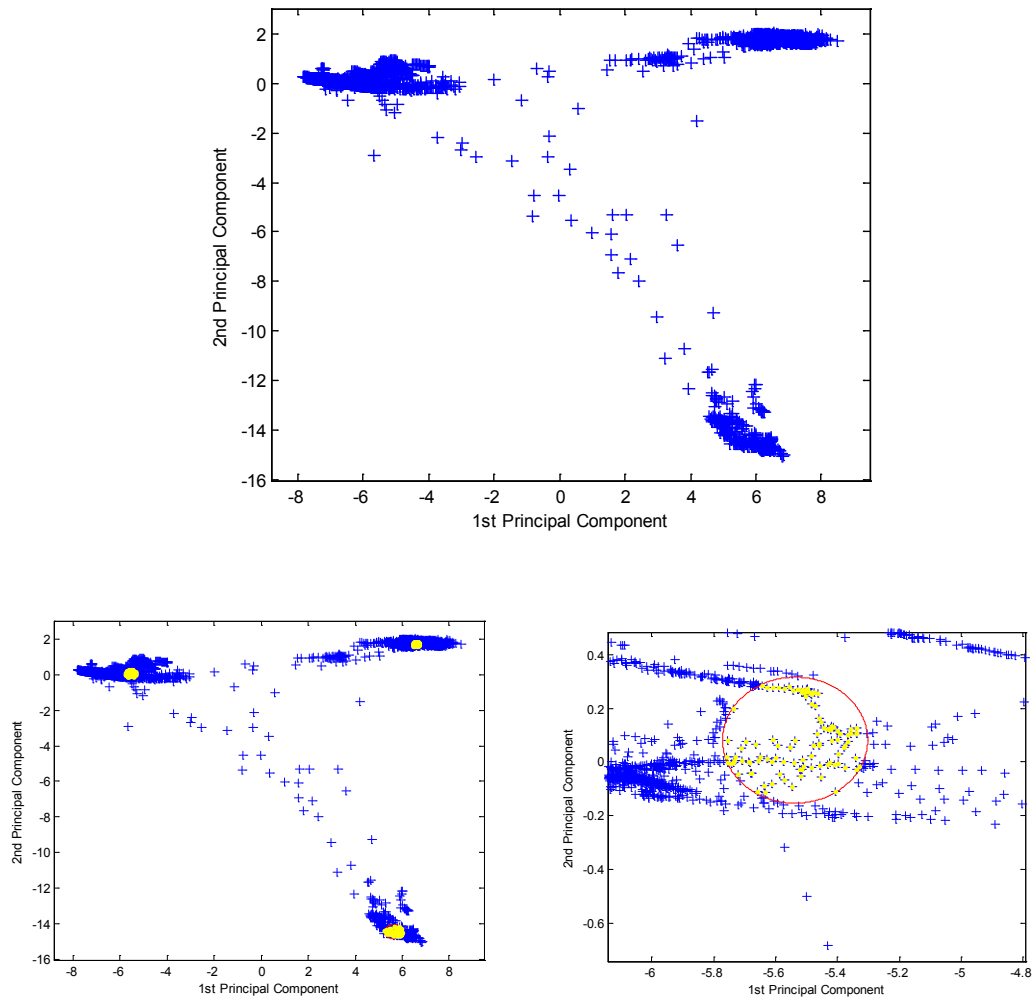


Figura 3.4: Observaciones de la base de datos de entrenamiento graficados en el espacio transformado vía PCA y selección de aquellas observaciones representativas que constituyen C1.

3.2.2.1.4 DIVISIÓN CONJUNTO INICIAL C1 (D)

Este bloque tiene como finalidad fraccionar el conjunto C1 en 4 subconjuntos que serán utilizados en la etapa de optimización del algoritmo principal. Esta división se realiza considerando el estadístico de Hotelling, T^2 , obtenido al analizar la matriz de errores generada por el modelo SBM inicial. El estadístico de Hotelling permite ordenar de forma descendente todas las observaciones contenidas en C1 y así, construir 4 subconjuntos en donde el primero de ellos es generado a partir de las observaciones de C1 que fueron peor estimadas según Hotelling; el segundo es generado con el segundo grupo de observaciones mal estimadas; y siguiendo el

criterio, el último subconjunto es constituido por observaciones en C1 que fueron mejor estimadas por el modelo SBM inicial, según Hotelling.

3.2.2.1.5 SBM (E)

El bloque “SBM” tiene como objetivo generar un modelo SBM en base a observaciones de entrada y salida predefinidos (Di y Do , respectivamente). Además realiza una estimación de las observaciones de salida de acuerdo a valores de las variables de entrada (Min). El procedimiento de generación del modelo y de las estimaciones realizadas sigue el descrito en Sección 2.3.2.2.

Como salida se tienen las estimaciones realizadas por el modelo para los instantes de tiempo asociados a Min ; un vector de pesos máximos de cada observación estimada por el modelo, la cual está ligada a la similitud de una observación con alguna presente en la matriz Di ; y por último, la media y desviación estándar de la matriz Do , utilizada para normalizar el error de estimación del modelo SBM.

3.2.2.1.6 SIMILITUD (F)

El bloque “Similitud” tiene como propósito determinar un valor escalar al grado de similitud de 2 vectores. Esta función de similitud puede ser la función triangular explicada en Sección 2.3.2.2. Tiene como entrada 2 vectores y una constante asociada al tamaño de los vectores. Como salida se encuentra un número escalar que representa el grado de similitud de ambos vectores, siendo h el valor para dos vectores idénticos y cero para 2 totalmente distintos.

3.2.2.1.7 VECTORES LD (G)

Este bloque tiene como objetivo principal eliminar filas linealmente dependientes o que son, en forma aproximada, linealmente dependiente. Considerándose aquellas que cumplen lo siguiente:

$$\det(v_i, v_j) \cong 0 \tag{3.1}$$

La creación de esta función se realiza debido a los constantes problemas que surgen cuando se necesita invertir la matriz relacionada con Di (dentro del proceso de cálculo del modelo SBM), en donde existe la posibilidad de encontrar vectores linealmente dependientes que imposibilitan al

computador el cálculo de inversión de forma correcta. Este bloque soluciona tal problema y es explicado a continuación.

Sea v_1, v_2, \dots, v_n las filas de la matriz de entrada Di , el algoritmo crea un ordenamiento descendiente de estos vectores utilizando como criterio la función de similitud explicada 3.2.2.1.6, es así como, cuando se tienen dos vectores muy similares, el parámetro de ordenamiento es máximo (h) y en caso contrario, el parámetro es cero. De esta manera es posible crear un ranking de vectores, considerando las similitudes entre cada uno de los vectores con el resto de ellos. Esto permite identificar las combinaciones más similares entre 2 vectores y eliminar uno de ellos (el más similar a los otros $(n-2)$ vectores). Este proceso se realiza de tal manera de eliminar un número definido de vectores (% del total de observaciones de Di), siendo estos los primeros de tal ranking.

En resumen, realizando un estudio de la similitud de cada uno de los vectores con los demás, se podrá escoger inteligentemente aquellos que tienen mayor similitud con el resto y así, reducir la dimensión de la matriz Di y el conjunto C1 que a su vez producirá una matriz con mayores posibilidades de ser invertida correctamente.

3.2.2.1.8 ERROR DE ESTIMACIÓN (H)

El bloque “Errores de estimación” tiene como propósito generar una matriz de error de estimación basada en la resta de la matriz de mediciones obtenidas en planta y de la matriz de estimaciones calculadas por un modelo SBM, es así como:

$$MError = MoutMedido - MoutEstimado \quad (3.2)$$

$MError$ puede ser suavizado utilizando el bloque suavizar. Además, se genera un vector que representa el error medio cuadrático de cada una de las observaciones estimadas. Donde para una secuencia de L entradas:

$$[x_i]_{i=1..L}, x_i \in R^{nV_{in}}, nV_{in} \text{ número de variables de entrada}$$

Su correspondiente secuencia de salidas

$$[y_i]_{i=1..L}, y_i \in R^{nV_{out}}, nV_{out} \text{ número de variables de salida}$$

Es estimada mediante: $[\hat{y}_i]_{i=1..L}, \hat{y}_i \in R^{n_{Vout}}$. Entonces, el error de estimación asociado estará definido por $[e_i]_{i=1..L}$, dónde:

$$e_i = \left(\frac{1}{n_{Vout}} \right) \|y_i - \hat{y}_i\|_2 \in R \quad \forall i = 1 \dots L \quad (3.3)$$

3.2.2.1.9 FILTRO MEDIANA MÓVIL (I)

El bloque “Filtro mediana móvil” tiene como propósito filtrar el error de estimación de cada una de las variables de salida encontradas en las columnas de la matriz de error de estimación. Para ello, es utilizada la mediana móvil con una ventana fija. La mediana móvil es descrita a continuación:

Sea $V \in R^L$, se define el vector filtrado vía mediana móvil como:

$$V_f(i) = \text{mediana}(\{V(i - \text{rangoVentana}), V(i - \text{rangoVentana} + 1), \dots, V(i), \dots V(i + \text{rangoVentana})\}) \quad (3.4)$$

$$\forall i \in [1 + \text{rangoVentana} \dots L - \text{rangoVentana}]$$

Y

$$V_f(j) = V(j) \quad (3.5)$$

$$\forall j \in [1 \dots \text{rangoVentana}] \cup [L - \text{RangoVentana} \dots l]$$

3.2.2.1.10 ANÁLISIS MEDIANTE HOTELLING (J)

El bloque “Análisis mediante Hotelling” tiene como finalidad obtener el estadístico de Hotelling de la matriz de errores normalizada, procedimiento descrito en 2.3.3.2 - Test estadístico de Hotelling. Luego, utiliza el bloque Selección vía Hotelling para seleccionar instantes de tiempo que debiesen estar en el conjunto de observaciones final que definirá el modelo SBM.

3.2.2.1.11 SELECCIÓN VÍA HOTELLING (K)

Este bloque realiza una selección de observaciones candidatas para ser consideradas en el conjunto que define un modelo SBM (C_{SBM}). Para ello, utiliza el estadístico de Hotelling, T^2 , generado por el bloque J y determina los instantes de tiempo en las cuales existe una mala

estimación, según un modelo SBM anterior. La forma de agregar los nuevos instantes es realizado dividiendo el vector de instantes T^2 en sub-ventanas (SV_i) de un largo igual a un porcentaje del largo total de T^2 . En cada una de estas ventanas se calcula el promedio del estadístico de Hotelling que se encuentren bajo el umbral deseado (según umbral de Hotelling para base de entrenamiento, ver Sección 2.3.3.2). De cada ventana se agregan los instantes que más cerca, por debajo, estén de tal umbral. La cantidad de observaciones a agregar desde cada una de la sub-ventanas es directamente proporcional al promedio calculado.

$$nOBSxVentana = nOBSAdicionales * \left(\frac{\text{promedio}(SV_i)}{\text{promedio}(T2)} \right) \quad (3.6)$$

3.2.2.1.12 AJUSTE FINAL (L)

Este bloque tiene como propósito ajustar el modelo SBM final añadiendo algunas observaciones que no fueron consideradas en las etapas anteriores. Esta selección final de elementos es opcional y será ejecutada si el usuario lo requiere. La intención es agregar observaciones de la base de datos de entrenamiento que tienen un error medio cuadrático fuera de lo normal y cuyo estadístico de Hotelling supera el umbral establecido en el último análisis vía Hotelling. Así, se agregan para mejorar significativamente el error del modelo SBM, obtenido como sigue:

$$E = \frac{1}{N} \sum_{i=1}^N e_i \quad (3.7)$$

N número de observaciones.

Y e_i calculado según Ecuación (3.3).

El número de observaciones adicionales es escogido por el usuario, además, para evitar agregar observaciones muy similares, el algoritmo excluye observaciones que se localizan a menos de un número definido de elementos de uno ya agregado.

3.2.2.2 BLOQUE PRINCIPAL: MODELO SBM

El bloque principal Modelo SBM tiene como objetivo seleccionar las observaciones que serán representativas del proceso en estudio. Recibe como entrada las matrices de entrada y salida de entrenamiento; el porcentaje de observaciones que constituirán el modelo SBM y que serán representativos del proceso; un índice inicial y final para la ubicación de elementos de la base de

datos de entrenamiento; y una variable booleana que indica si se desea graficar los resultados de cada una de las etapas que constituyen este algoritmo. Estas etapas son ilustradas en la Figura 3.3.

La primera fase del algoritmo principal tiene por objetivo seleccionar un conjunto inicial de observaciones C1 para así generar un modelo SBM inicial, esto es hecho con el bloque A (según Tabla 3.1). El modelo SBM inicial es generado utilizando los bloques E, F, y G.

La estimación de las salidas obtenidas con el modelo inicial es utilizada para realizar un tratamiento mediante análisis de componentes principales y test de Hotelling a la matriz de errores, utilizando para ello los bloques H, I y J. De esta manera y a través del bloque K, es posible incluir observaciones no consideradas en el conjunto inicial (4% de la base de datos) y así, completar un conjunto preliminar C2, constituido por un 10 % de las observaciones de la base de entrenamiento. Posteriormente se utilizan, nuevamente, los bloques E,F,G,H,I y J para generar un modelo SBM preliminar; obtener la matriz de errores de estimación y su respectivo vector estadístico de Hotelling. Con este vector es posible seleccionar observaciones candidatas a través del bloque K, el cual es utilizado en el proceso de optimización del modelo SBM.

La etapa de optimización está constituida por 4 sub-etapas que consideran la eliminación parcial del conjunto inicial de observaciones (ver Figura 3.5). Esta eliminación se hace debido a que existen observaciones dentro de tal conjunto que pueden ser redundantes (entregan misma información al modelo SBM) y/o no son representativas de las verdaderas condiciones de operación del proceso. Es por ello que, se modifica el conjunto C2 de observaciones que definen el modelo preliminar, reemplazando un 1% de los instantes de tiempo seleccionados inicialmente (encontrados en C1) por observaciones candidatas obtenidas por el último modelo SBM creado y así, se crea un conjunto C3 de observaciones, para luego evaluar la calidad del nuevo modelo SBM creado con tal conjunto, conservando el 1% de observaciones candidatas si es que el modelo mejora en la estimación de las salidas, y en caso contrario, se reincorpora el 1% de las observaciones que fueron reemplazadas. Siguiendo este procedimiento, se realizan 4 iteraciones que permiten ajustar el conjunto de observaciones C3 que definirán el modelo SBM final.

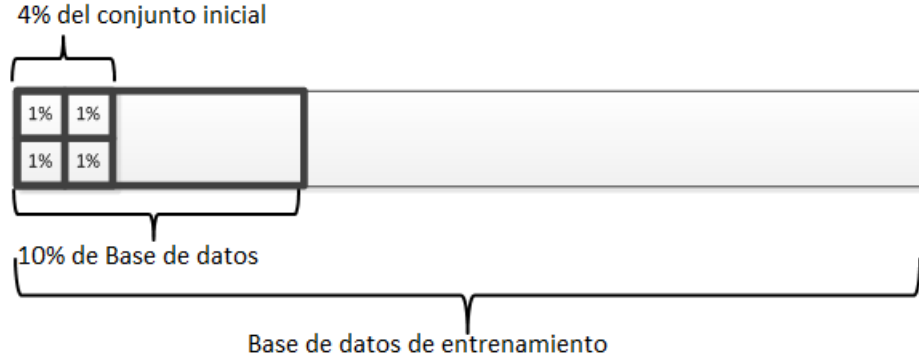


Figura 3.5: Subdivisión del conjunto inicial C1.

Al finalizar el proceso de optimización, el algoritmo principal ofrece al usuario ajustar el modelo SBM final usando L. Este ajuste consiste en adicionar observaciones para mejorar la calidad de las estimaciones del modelo final con lo que el conjunto de observaciones seleccionadas (C_{SBM}) queda finalmente establecido.

La salida del algoritmo principal es un objeto estructura de MATLAB[®] la cual consta de los siguientes elementos.

- $D_i \in \mathbf{R}^{m \times nVin}$: Matriz con $nVin$ variables de entrada y m observaciones, siendo $m \cong 0.1 \cdot n$.
- $D_o \in \mathbf{R}^{m \times nVout}$: Matriz con $nVout$ variables de salida y m observaciones, siendo $m \cong 0.1 \cdot n$.
- $ETrain \in \mathbf{R}^{n \times nVout}$: Matriz que representa el error de estimación del modelo SBM final, para cada variable de salida y para cada una de las observaciones.
- $MeanETrain \in \mathbf{R}^{1 \times nVout}$: Vector que contiene en cada elemento el promedio aritmético de cada una de las columnas de $ETrain$.
- $SigmaETrain \in \mathbf{R}^{1 \times nVout}$: Vector que contiene en cada elemento la desviación estándar de cada una de las columnas de $ETrain$. En el caso en que una columna tenga desviación estándar igual a cero, se reemplaza tal valor por uno.
- $vpETrain \in \mathbf{R}^{nVout \times nVout}$: Matriz que contiene los vectores propios de la matriz de covarianza de $ETrain$. Representan cada uno de las componentes principales obtenidas al hacer PCA a la matriz $ETrain$.
- $vlETrain \in \mathbf{R}^{nVout \times 1}$: Matriz que contiene los valores propios de la matriz de covarianza de $ETrain$.

Los elementos mencionados anteriormente son utilizados en la herramienta de detección de anomalías en línea, en donde las matrices D_i y D_o son las que permiten generar un modelo SBM y así estimar, para un instante de tiempo en particular, los valores para las variables de salida dada las mediciones de las variables de entrada.

La matriz $ETrain$ es utilizada en la herramienta de detección de anomalías para comparar gráficamente las estimaciones realizadas en la base de datos de entrenamiento y los nuevos datos que se vayan incorporando.

Los vectores $MeanETrain$ y $SigmaETrain$ son utilizados para normalizar las nuevas salidas estimadas a través del modelo SBM manteniendo un rango de valores común, lo cual es importante para el análisis de los errores obtenidos basado en el test de Hotelling.

La matriz $vpETrain$ y el vector $vlETrain$ son necesarios para transformar los nuevos errores obtenidos por el modelo SBM al espacio transformado generado por la base de datos de entrenamiento, transformación requerida para obtener resultados coherentes en el cálculo del estadístico de Hotelling.

En la Figura 3.6 se presenta un diagrama de flujo que permite seguir paso a paso las acciones que toma el algoritmo principal para generar un modelo SBM final, se indican las funciones utilizadas en cada bloque a través de una letra que hace referencia a la ordenación indicada en la Tabla 3.1.

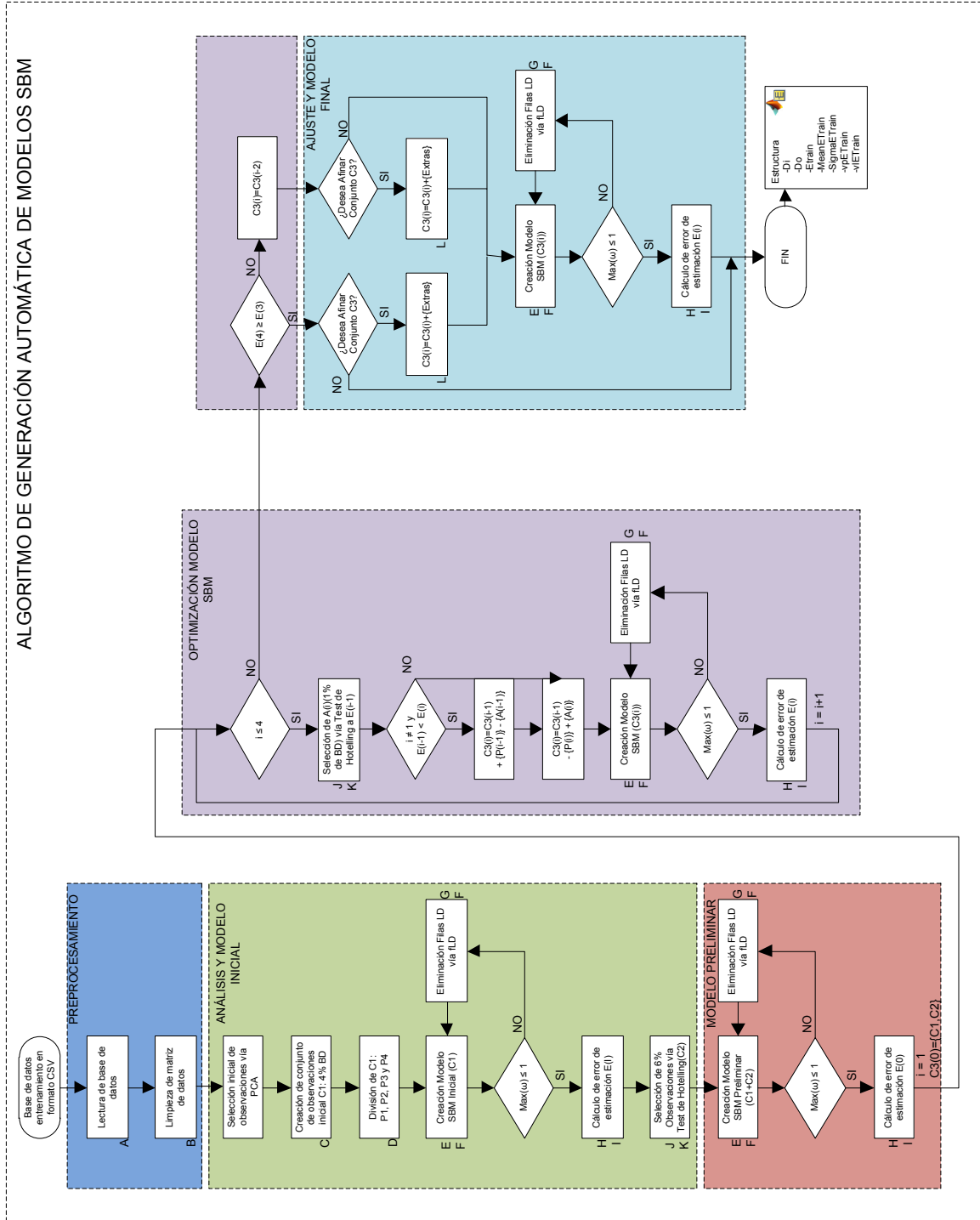


Figura 3.6: Diagrama de Flujo – Descripción detallada de algoritmo de generación automática de modelos SBM.

El algoritmo antes descrito permite generar modelos SBM de procesos industriales multivariados de forma automatizada. Esta herramienta será utilizada en el sistema de detección de anomalías

propuesta en este trabajo; sin embargo, es necesario diseñar una metodología que permita estudiar cual es la mejor selección de las variables de tal modelo, el cual concluirá en la generación de alarmas y detección de anomalías, es por ello que, se describe a continuación el procedimiento diseñado para identificar aquellas variables más significativas en la detección de una anomalía en particular.

3.3 METODOLOGÍA PARA ESTUDIO DE EVENTOS USANDO MODELOS SBM.

Luego de entregada una descripción completa del algoritmo de generación automática de modelos SBM, se presenta a continuación una metodología de estudio de eventos basado en la modelación SBM y el algoritmo creado en Sección 3.2.

Como se mencionó en tal etapa, la elección de variables tanto de entrada como de salida de un modelo SBM debe realizarse en base a un estudio preliminar que utiliza información de diversas fuentes y métodos. Es por ello que, se diseñó una metodología que permite seleccionar variables de entrada y salida con el motivo hacer más efectiva la herramienta de detección de anomalías vía modelos SBM. Definiendo una anomalía en particular como *evento* se presenta un diagrama de flujo que ilustra el procedimiento que debiese terminar con la elección de las variables de entrada y salida que producen mayor efectividad en la detección del mismo.

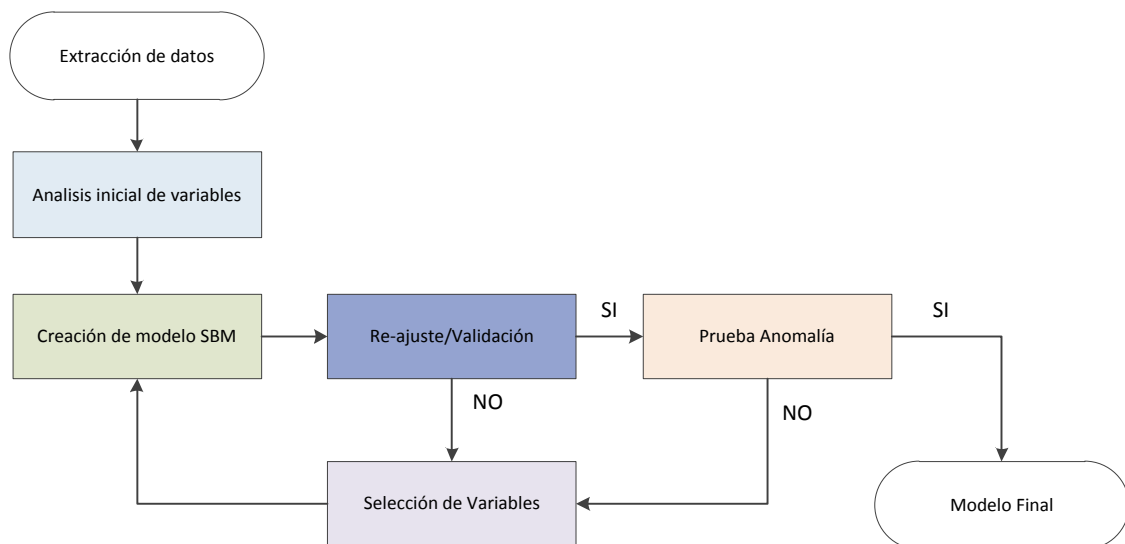


Figura 3.7: Diagrama de flujo de metodología para el estudio de eventos usando modelos SBM.

El diagrama de flujo anterior está compuesto por 7 etapas las cuales serán detalladas a continuación.

3.3.1 EXTRACCIÓN DE DATOS

Esta etapa cumple la misma función que la etapa pre-procesamiento vista en Sección 3.2 - Algoritmo de generación automática de modelos SBM. Es decir, el procedimiento se inicia con una base de datos de entrenamiento, la cual debe ser representativa del proceso en estudio, rica en información sobre condiciones de operación y lo suficientemente válida para su análisis. Es posible utilizar los mismos bloques vistos en Sección 3.2.2.1. Es así como, se tienen datos válidos y limpios sobre la operación del proceso que considera todas las variables disponibles en la base de datos de entrenamiento.

3.3.2 ANÁLISIS INICIAL DE VARIABLES

En esta etapa se determinan los conjuntos de variables de entrada y de salida iniciales. Se considera una etapa crucial para el éxito de la herramienta propuesta en este trabajo, pues es aquí donde se elimina la mayor parte de las variables que no aportan a la solución del problema y que introducen mayoritariamente ruido.

La determinación de variables se realiza en base a un estudio previo del proceso, la cual involucra información proveniente de operadores y especialistas. Además se realiza una verificación de las variables para eliminar aquellas redundantes en su información o repetidas.

Adicionalmente, esta etapa considera un estudio basado en el método PLS, básicamente para la eliminación de variables. Basado en lo mencionado en Sección 2.3.2.1 y Sección 3.1, el análisis PLS permite seleccionar/eliminar variables, tanto de entrada como de salida, determinando cuales de ellas son las que introducen ruido y no están aportando a la descripción de la estructura de los datos del proceso.

El procedimiento propuesto empieza al escoger un conjunto inicial de variables de entrada y salida, para luego analizar los datos de tales variables con regresión PLS y así, graficar los vectores de carga verificando la ubicación de cada una de las variables dentro de tal gráfico mediante el cálculo de su distancia al origen. De esta forma, las variables más cercanas a tal punto son exactamente las variables que no son explicadas por los vectores de carga y, por lo

tanto, no aportan a la correcta estimación del modelo. Finalmente, es posible eliminar variables de entrada y de salida, escogiendo un umbral apropiado en donde todas las variables cuya distancia al origen sea menor al umbral quedan eliminadas del modelo SBM.

3.3.3 CREACIÓN DE MODELO

Luego de escogido un conjunto inicial de variables de entrada y salida en base a lo explicado en el punto anterior, se está en condiciones de generar el primer modelo SBM. Esta generación se realiza utilizando el algoritmo descrito en Sección 3.2.

Es posible que etapas posteriores reajusten el conjunto de variables de entrada y salida para el modelo y, en tal caso, la metodología requerirá de esta etapa para generar un nuevo modelo SBM.

3.3.4 RE-AJUSTE/VALIDACIÓN

En esta etapa se realiza una validación del modelo obtenido en base a la elección de variables de entrada y salida de Sección 3.3.2. Esto requiere utilizar una base de datos de validación que corresponda a las condiciones de operación que debiesen ser bien modeladas. Además, en esta etapa se realiza un ajuste del umbral de Hotelling en caso de que existan observaciones de la base de datos de validación que sobrepasan el umbral definido en la Ecuación (2.45). Como última acción, se realiza una revisión de variables que son significativamente mal estimadas, siendo necesario eliminarlas y/o considerarlas como variables de entrada, esta tarea es realizada en la etapa “Selección de variables”

3.3.5 PRUEBA ANOMALÍA

Esta fase de la metodología de estudio de eventos está relacionada con la prueba del modelo SBM utilizando una base de datos que contiene información del evento en estudio. Esta etapa permite comprobar la efectividad del modelo SBM en la determinación del evento. La detección del evento en cuestión se realiza por medio del umbral de Hotelling de Ecuación (2.46) y ajustado en Sección 3.3.4. En casos en que el modelo SBM no cumpla con el objetivo de detección, se procede a observar las variables estimadas, verificando cuales de ellas presentan problemas y así, se redefine el conjunto de variables de entrada y salida en la etapa “Selección de variables”.

3.3.6 SELECCIÓN DE VARIABLES

En base a los resultados preliminares obtenidos en Sección 3.3.4 y Sección 3.3.5 se eliminan/reubican variables que no aportan a la detección. Luego de esto, se procede a generar un nuevo modelo SBM según el conjunto de variables de entrada y salida actualizado.

3.3.7 MODELO FINAL

Realizadas las etapas anteriores, se está en condiciones de construir un modelo SBM final, que considera las variables que realmente aportan a la determinación del evento en estudio.

3.4 DETECCIÓN DE ANOMALÍAS EN PROCESOS INDUSTRIALES USANDO MODELOS BASADOS EN SIMILITUD

Luego de descrito el algoritmo de generación automática de modelos SBM y la metodología de estudio de eventos, se procede a describir la herramienta que hará uso de tales procedimientos para detectar anomalías en procesos industriales multivariados en línea.

La herramienta de detección de anomalías usando modelos SBM sigue una serie de etapas las cuales son ilustradas en el siguiente diagrama de flujo.

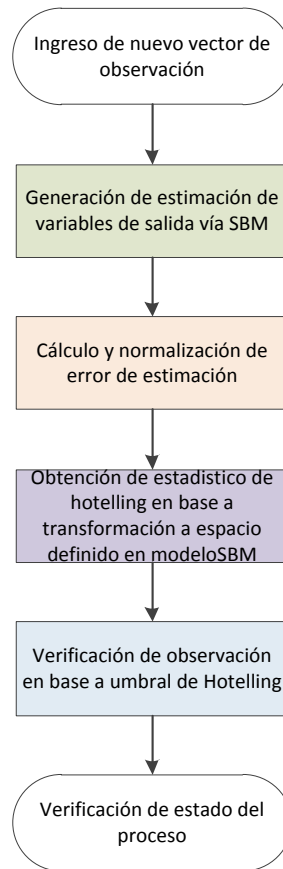


Figura 3.8: Diagrama de flujo de herramienta de detección de anomalías usando modelos SBM.

Como se puede observar en la Figura 3.8, la herramienta de detección se ejecuta cada vez que una nueva observación es obtenida desde el proceso. Luego de ello, se generan las estimaciones del modelo SBM utilizando para ello las matrices D_i y D_o (características de un modelo SBM definido). Realizadas las estimaciones se procede a compararlas con las variables medidas para generar un vector de error de estimación el cual es normalizado usando como media $mediaETrain$ y como desviación estándar $sigmaETrain$. El error de estimación normalizado es posteriormente transformado al espacio caracterizado por $vLETrain$ y $vpETrain$. Finalmente, se obtiene el estadístico de Hotelling de tal vector de observación el cual es comparado con el umbral de Hotelling previamente definido. De esta manera es posible definir para cada vector de observación un estado del proceso. Cuando el estadístico de Hotelling para la observación k esté bajo el umbral se dirá que el proceso está operando en condición de operación normal, y en caso contrario, cuando el estadístico de Hotelling se encuentre sobre el umbral se dirá que el proceso se encuentra operando anormalmente.

El sistema de detección de anomalías en procesos industriales usando modelos SBM propuesto necesita de una correcta evaluación, en la que se verifique tanto la capacidad de modelamiento de los modelos SBM como también, la capacidad de detección de eventos anómalos. Esta evaluación es presentada en el capítulo siguiente, el cual entrega los resultados obtenidos al utilizar dicho sistema en un proceso industrial existente y con anomalías conocidas puntuales.

CAPÍTULO 4. ESTUDIO DE EVENTO – RESULTADOS

El proceso de detección de anomalías basado en residuos está dividido en dos etapas fundamentales: *Generación de residuos y análisis de residuos* [35]. En la generación de residuos basada en observadores se requiere un modelo que permita emular el comportamiento del proceso y de esta forma compararlo con el comportamiento medido del mismo. El análisis de residuos es la etapa en la cual se estudian los residuos, con el propósito de extraer información de ellos y evaluarla. Propuestos los métodos de generación y análisis de residuos a utilizar en este trabajo, solo queda generar resultados para su análisis posterior.

En este capítulo se procede a entregar los resultados de la aplicación de las herramientas generadas anteriormente con el fin de evaluar su rendimiento y capacidad de detección. Esta aplicación necesita de un proceso industrial multivariado, con información disponible de operación normal como también, de condiciones anómalas identificadas. Por este motivo, la Sección 4.1 entrega una descripción de un proceso industrial multivariado y de las bases de datos históricas que serán utilizadas para evaluar la herramienta de detección de anomalías usando modelos SBM. La Sección 4.2 entrega los resultados obtenidos al evaluar la capacidad de modelamiento del algoritmo de generación automática de modelos SBM (visto en Sección 3.2), esta capacidad será comparada con respecto al rendimiento obtenido al utilizar modelos estáticos y lineales en los parámetros. Finalmente, la Sección 4.3 entrega los resultados conseguidos, al utilizar modelos SBM (construidos usando el algoritmo de generación automática) y la metodología de estudio de eventos (visto en Sección 3.3), en la detección en línea de anomalías en procesos industriales.

4.1 DESCRIPCIÓN DEL PROCESO INDUSTRIAL ESTUDIADO Y DE LAS BASES DE DATOS UTILIZADAS

Como ya ha sido mencionado, en la actualidad existe una gran diversidad de industrias que están constituidas por procesos que incorporan rutinas de monitoreo y algoritmos de detección de eventos anómalos, una de ellas es la industria de la generación de electricidad. Estas industrias necesitan de sistemas de protección y de detección oportuna de anomalías debido a los altos costos económicos asociados a la detención del proceso y los efectos negativos que pueda provocar en la red de generación en la cual está conectada. Uno de los sistemas más complejos de analizar es la generación termoeléctrica de ciclo combinado, este tipo de generación constituye el 20.0 % de la producción nacional de electricidad del año 2010 en Chile², y un 23.01 % de la producción de electricidad del año 2010 en España³. Parte de proceso industrial es analizado y estudiado con las herramientas creadas y presentadas en este trabajo, por lo que se hace necesario realizar una breve descripción de los componentes que lo constituyen y su funcionamiento.

4.1.1 CENTRAL TÉRMICA DE CICLO COMBINADO

Una central térmica de ciclo combinado es aquella donde se genera electricidad mediante la utilización conjunta de dos máquinas generadoras: Un turbogruppo de gas/diesel y un turbogruppo de vapor. Es decir, para la transformación de la energía del combustible en electricidad se superponen dos ciclos [32]:

- a) El ciclo de Brayton (turbina de gas/diesel): Toma el aire directamente de la atmósfera y se somete a un calentamiento y compresión para aprovecharlo como energía mecánica o eléctrica.
- b) El ciclo de Rankine (turbina de vapor): Donde se relaciona el consumo de calor con la producción de trabajo o creación de energía a partir de vapor.

² Según Informe Anual 2010, Energía Eléctrica, INE Chile.

³ Según Informe del sistema eléctrico del 2010, Red Eléctrica de España.

Las partes que conforman una central térmica de ciclo combinado son las siguientes (ver Figura 4.1):

- **Turbina de gas/diesel:** Consta de compresor, cámara de combustión y la propia turbina.
 - **Compresor:** Generalmente es un compresor por etapas y su función es inyectar el aire a presión por la combustión del gas y la refrigeración de las zonas calientes.
 - **Cámara de combustión:** en este punto de la instalación es donde se mezclan el gas natural o diesel con el aire a presión y se produce la combustión.
 - **Turbina de gas:** en ella se produce la expansión de gases que provienen de la cámara de combustión. Consta de tres o cuatro etapas de expansión y la temperatura de los gases en la entrada está alrededor de 1.400°C saliendo de la turbina a temperaturas superiores a los 600°C .
- **Caldera de recuperación:** En esta caldera convencional el calor de los gases que provienen de la turbina de gas se aprovechan en un ciclo de agua-vapor.
- **Turbina de vapor:** Esta turbina acostumbra a ser de tres cuerpos y está basada en la tecnología convencional. Es muy habitual que la turbina de gas y la turbina de vapor se encuentren acopladas a un mismo eje de manera que accionan un mismo generador eléctrico.

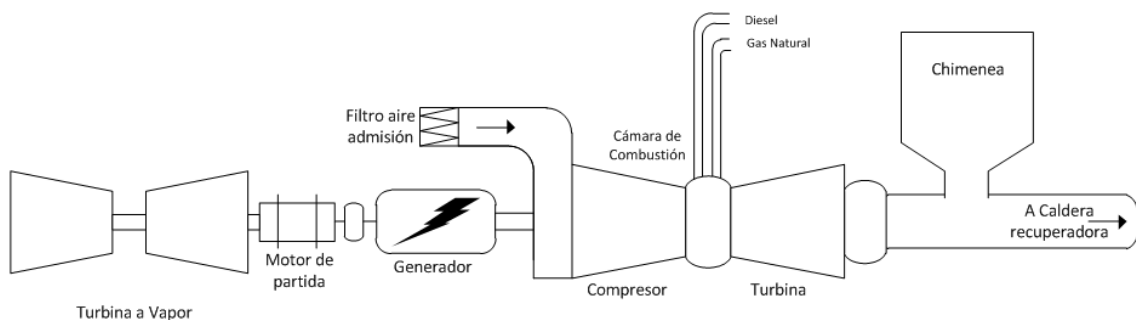


Figura 4.1: Esquema básico de Turbina a Gas/Diesel de una central térmica de ciclo combinado.

El funcionamiento de una central de este tipo sigue un procedimiento el cual se inicia cuando el aire es comprimido a alta presión en el compresor, después pasa a la cámara de combustión donde se mezcla con el combustible. A continuación, los gases de combustión pasan por la turbina de

gas donde se expansionan y su energía calorífica se transforma en energía mecánica transmitiéndolo al eje de la turbina.

Los gases que salen de la turbina de gas se llevan a una caldera de recuperación de calor para producir vapor, luego se inicia un ciclo agua-vapor utilizado en las centrales térmicas convencionales. A la salida de la turbina el vapor se condensa (transformándose nuevamente en agua) y vuelve a la caldera para empezar un nuevo ciclo de producción de vapor.

En particular, los algoritmos de detección de anomalías propuestos en este trabajo utilizarán la información de operación de las variables localizadas en la turbina a gas del sistema antes descrito (ver Figura 4.1), de esta forma se tendrán medidas de variables como temperaturas, presiones, vibraciones, entre otras; tanto en el compresor como en la cámara de combustión y turbina. Información que constituirán las bases de datos de entrenamiento, validación y prueba para los esquemas de detección desarrollados.

Como primera aproximación, el sistema de detección propuesto se enfocará al estudio de eventos anómalos de lenta evolución, como son desgastes y/o malfuncionamiento de componentes del proceso. De esta manera, se busca identificar prematuramente las condiciones que conllevan gradualmente a una operación anormal y así responder de forma temprana a tales eventos. Lo anterior esta basado en las bajas posibilidades de actuar en la prevención y detección de una anomalía abrupta de un proceso industrial de gran envergadura, sin embargo, el sistema de detección propuesto está capacitado a responder frente a comportamientos anómalos previos y recurrentes y que, eventualmente, terminen en una anomalía abrupta.

4.1.2 BASES DE DATOS UTILIZADAS

Las bases de datos utilizadas en la elaboración del sistema de detección de anomalías fueron proporcionadas por ENDESA-CHILE. Básicamente, se trata de 2 grandes archivos recopilatorios que contienen información sobre el funcionamiento de dos turbinas de Gas/Diesel. Las cuales son especificadas a continuación.

	Base de Datos T1	Base de Datos T2
N° Variables	192	203
N° Observaciones	7000 (6919 Válidas)	19579 (19530 Válidas)
Tiempo de Muestreo	10[<i>min</i>]	10[<i>min</i>]
N° condiciones de operación conocidas	3	3
N° de anomalías identificadas	0	1

Tabla 4.1: Bases de datos utilizadas.

La base de datos T1 contiene información de operación de una turbina T1 sin anomalías identificables, por lo que fue utilizada en la etapa de creación del algoritmo de generación automática de modelos SBM detallada en Sección 3.2. Esta base contiene datos de la turbina en operación y no operación (ver Figura 4.2 y Figura 4.3). Además, la turbina presenta dos condiciones de operación disímiles, operación con combustible GAS y operación con combustible DIESEL (ver Figura 4.4).

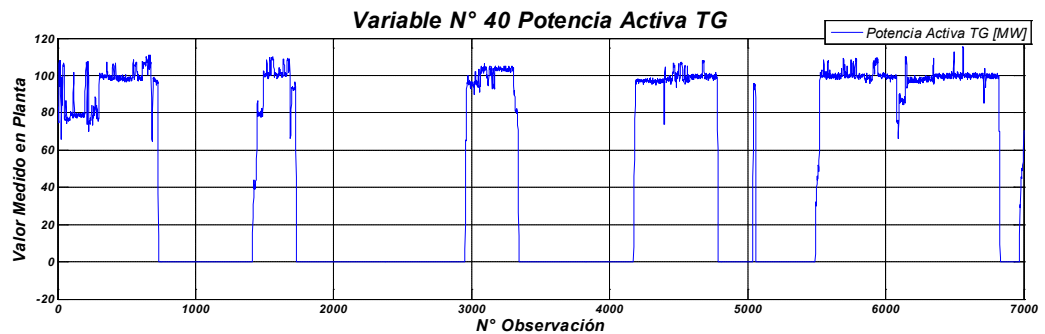


Figura 4.2: Potencia activa T1 [MW].

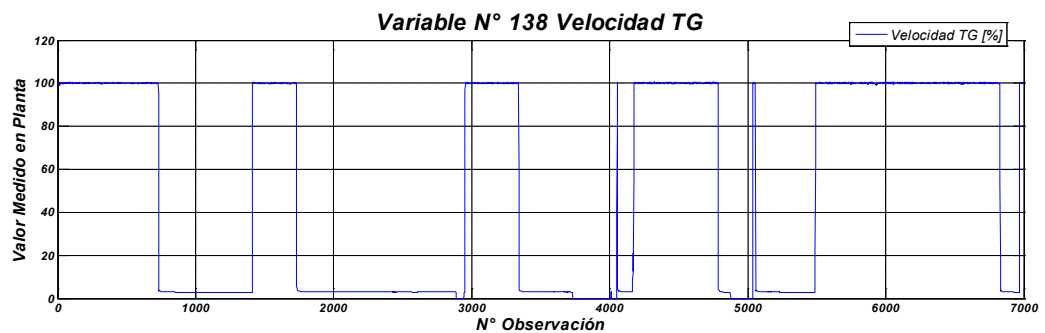


Figura 4.3: Velocidad turbina T1 [%].

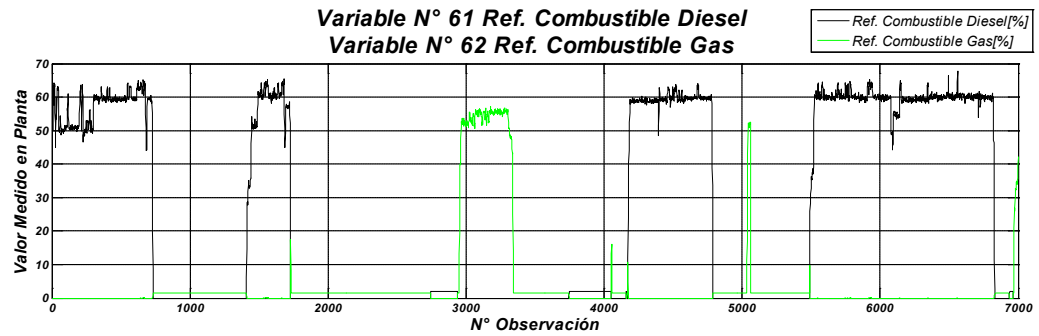


Figura 4.4: Referencia señal de combustible T1, DIESEL (en negro) y GAS (en verde).

La base de datos T2 contiene información de operación de una turbina T2 con una anomalía o evento identificado por operadores/especialistas, siendo ésta la base de datos utilizada para probar la metodología de estudio de eventos y el sistema de detección de anomalías presentados en 3.3 y 3.4, respectivamente. Esta base contiene datos de la turbina en operación a GAS y en no operación. Condiciones ilustradas en la Figura 4.5, Figura 4.6 y Figura 4.7.

La información presentada, referida a las bases de datos utilizadas, está limitada debido a los contratos de confidencialidad explicados en la Sección 1.4 - Indicación sobre confidencialidad .

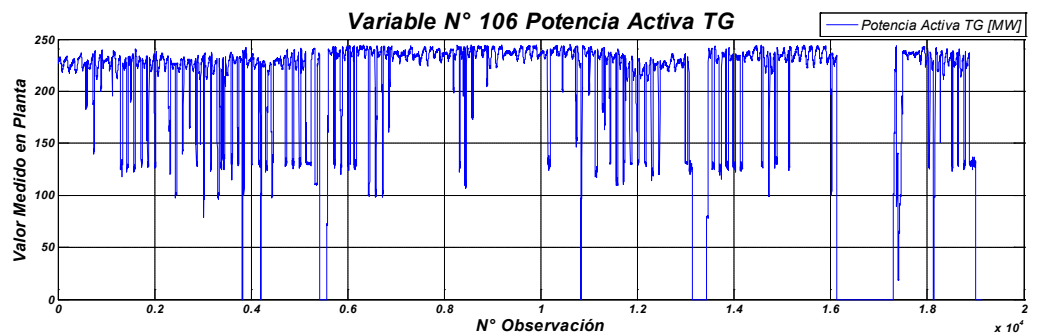


Figura 4.5: Potencia activa T2 [MW].

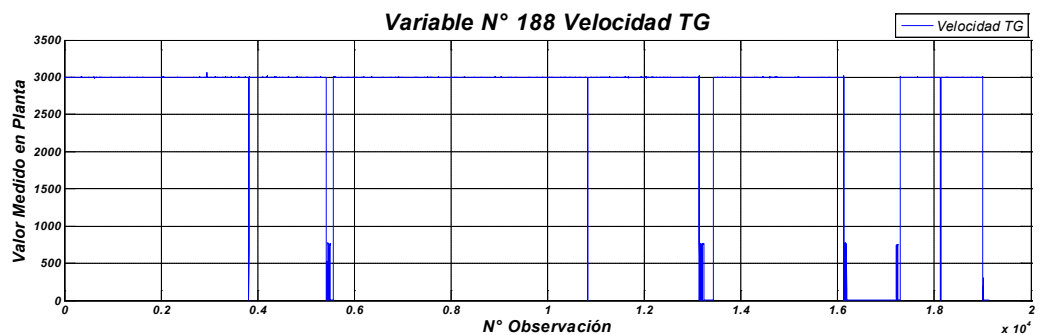


Figura 4.6: Velocidad turbina T2.

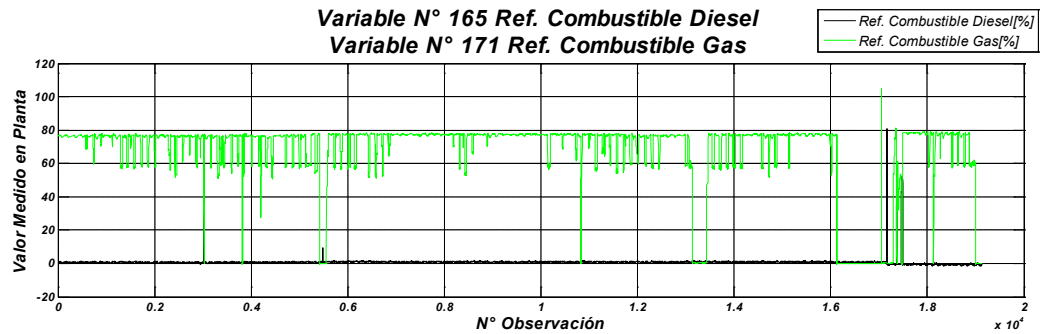


Figura 4.7: Referencia señal de combustible T2, DIESEL (en negro) y GAS (en verde).

Los gráficos antes mostrados permiten tener una visión general del comportamiento de las turbinas durante el estudio de las mismas. Se recalca la operación de las turbinas en distintos puntos de operación así como también la inclusión de datos de ambas turbinas detenidas. A continuación, se entregan los resultados obtenidos por el sistema de generación automática de modelos SBM sobre la base de datos T1 y aquellos resultados obtenidos por el sistema de detección diseñado en la base de datos T2.

4.2 RESULTADOS OBTENIDOS EN PRUEBA DE ALGORITMO DE GENERACIÓN DE MODELOS SBM Y BASE DE DATOS T1.

Luego de entregada una descripción general del proceso industrial a estudiar, se procede a entregar resultados de las distintas etapas que constituyen el algoritmo de generación automática de modelos SBM, con el fin de evaluar la capacidad de modelamiento de esta herramienta en contraste con otros tipos de modelos, en particular se hará una comparación con modelos lineales en los parámetros.

Para el análisis del algoritmo de generación automática de modelos SBM descrito en Sección 3.2 se utiliza la base de datos T1, pues esta no contiene anomalías conocidas ni identificadas. Además, presenta 3 condiciones de operación bien definidas: Operación con combustible gas, operación con combustible diesel y no operación. Dado que la importancia de esta etapa es evaluar la evolución de los modelos SBM implementados a lo largo de todo el algoritmo y de su capacidad de modelación, se utilizará un conjunto de variables de entrada y salida escogido sólo en base al conocimiento previo acerca del proceso; considera variables como presión, temperatura

y flujos en los distintos componentes que constituyen la turbina. De esta manera, dadas 192 variables se escogen los siguientes subconjuntos⁴:

Variables de Entrada N°	51	54	55	57	58	59	60	61
	64	71	75	76	80	81	96	97
	130	131	140	190				
Variables de Salida N°	2	3	8	9	10	11	12	13
	14	15	16	19	20	31	35	36
	37	42	43	44	45	46	47	48
	49	50	63	83	86	99	100	101
	103	104	120	122	125	127	129	134
	135	137	141	184	185	186	191	

Tabla 4.2: Conjunto inicial de variables de entrada y salida.

En primer lugar, se presentan los resultados obtenidos en la etapa de análisis inicial del algoritmo. Como se describe en Sección 3.2.1, el primer gráfico que entrega el algoritmo representa las observaciones en el espacio transformado por PCA. Como se puede apreciar en la Figura 4.8, es posible visualizar 3 agrupaciones de datos, los cuales representan los 3 puntos de operación correspondientes a la base de datos T1.

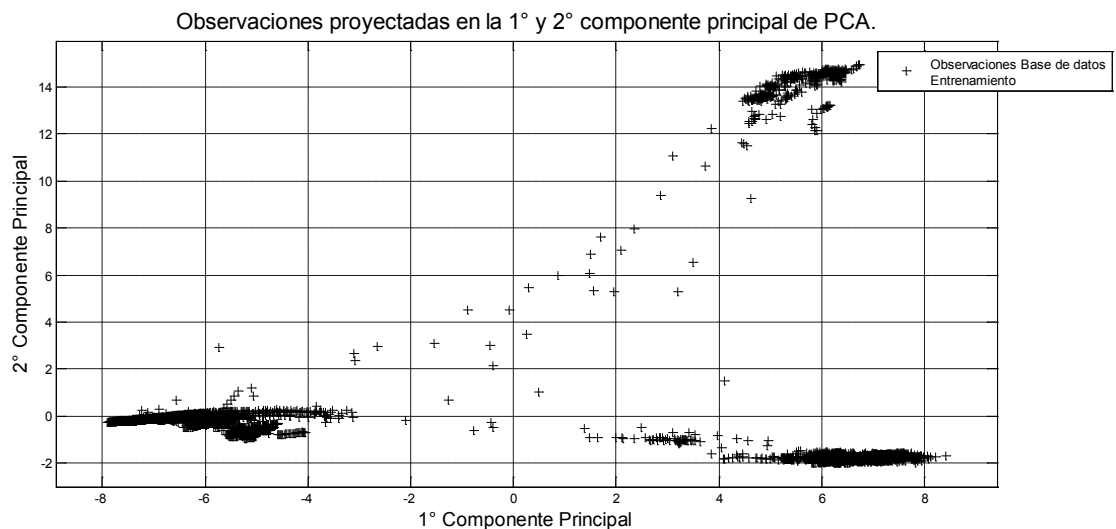


Figura 4.8: Observaciones proyectadas en la 1° y 2° componente principal de PCA.

⁴ Las variables de entrada y salida utilizadas son identificadas sólo por su índice de la base de datos, de esta manera se respetan los contratos de confidencialidad del proyecto en el cual se enmarca este trabajo, detallado en Sección 1.4.

La identificación de estas 3 agrupaciones puede ser realizada de forma manual, en tal caso el usuario deberá introducir los centros aproximados de cada una de las agrupaciones o, de forma automática, cuya selección de centros es realizada por el algoritmo, el cual asigna a cada una de las observaciones a un determinado grupo. Al realizar la selección de centros de manera automatizada se tiene una clasificación como la mostrada en la Figura 4.9, en donde se distingue cada observación a través de un color que identifica el modelo de operación asignado. Para validar esta selección, la Figura 4.10 muestra la potencia activa (Figura 4.10a) y la señal de combustible (Figura 4.10b) a lo largo de todas las observaciones registradas en T1. Esto permite identificar el modo de operación en que se encontraba el proceso para cada observación y que, al ser comparado con la clasificación hecha por el algoritmo (Figura 4.10c), es posible comprobar con facilidad que tal ordenación corresponde exactamente a los modos de operación registrados. Es así como; el modo de operación 1 representa la operación de la turbina con combustible diesel (en azul); el modo de operación 2 representa la operación de la turbina con combustible gas (en verde) y; el modo de operación 3 representa la no operación de la turbina (en rojo).

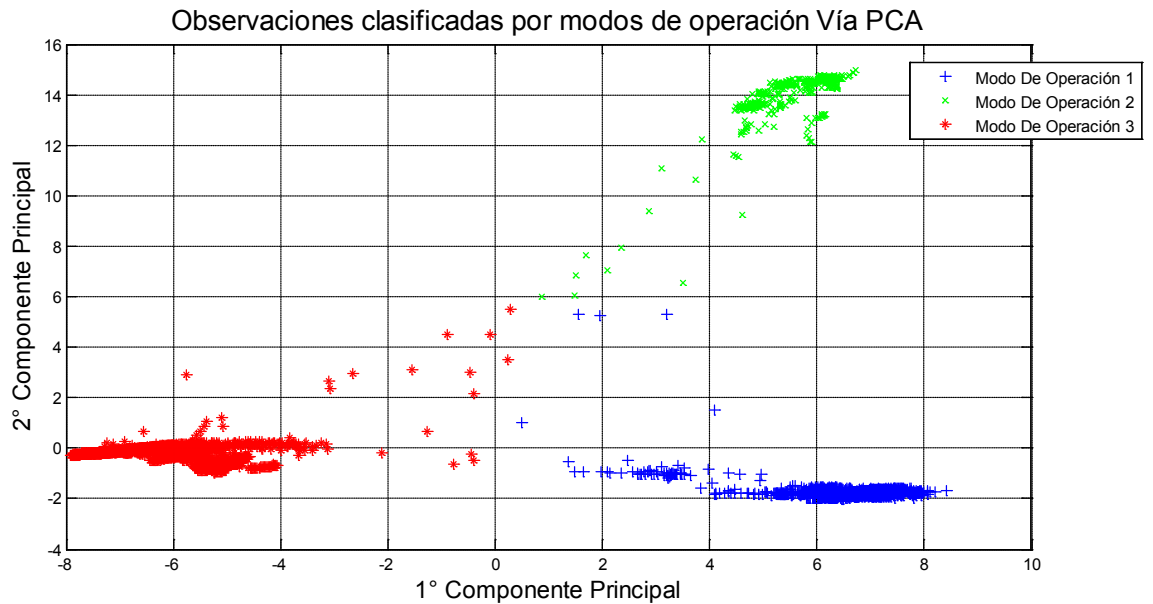


Figura 4.9: Observaciones clasificadas por modos de operación vía PCA.

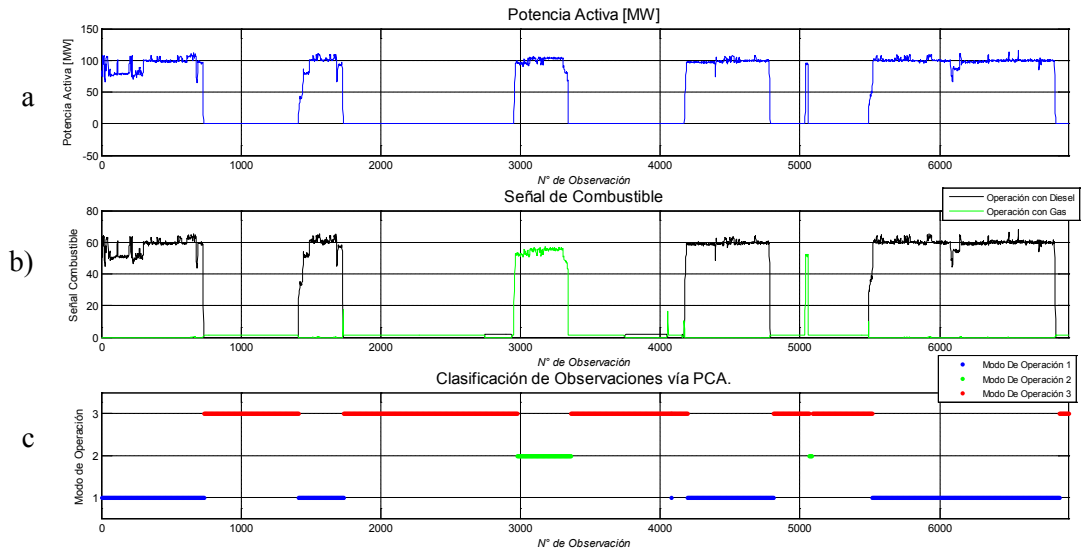


Figura 4.10: Comparación entre condiciones de operación identificables en a) y b), y la clasificación automática para cada observación c).

Estos resultados permiten concluir de forma preliminar que la metodología utilizada es exitosa, pues se logra identificar los puntos de operación del proceso y los centros de cada una de las agrupaciones. De esta forma, determinados los centros de las agrupaciones de forma manual o automática, el algoritmo procede a seleccionar un 4 % de las observaciones escogiendo umbrales elípticos en torno a cada uno de tales centros, como se observa en la Figura 4.11, y así incorporar al conjunto inicial (C_1) a todas las observaciones que queden localizadas dentro de cada elipse.

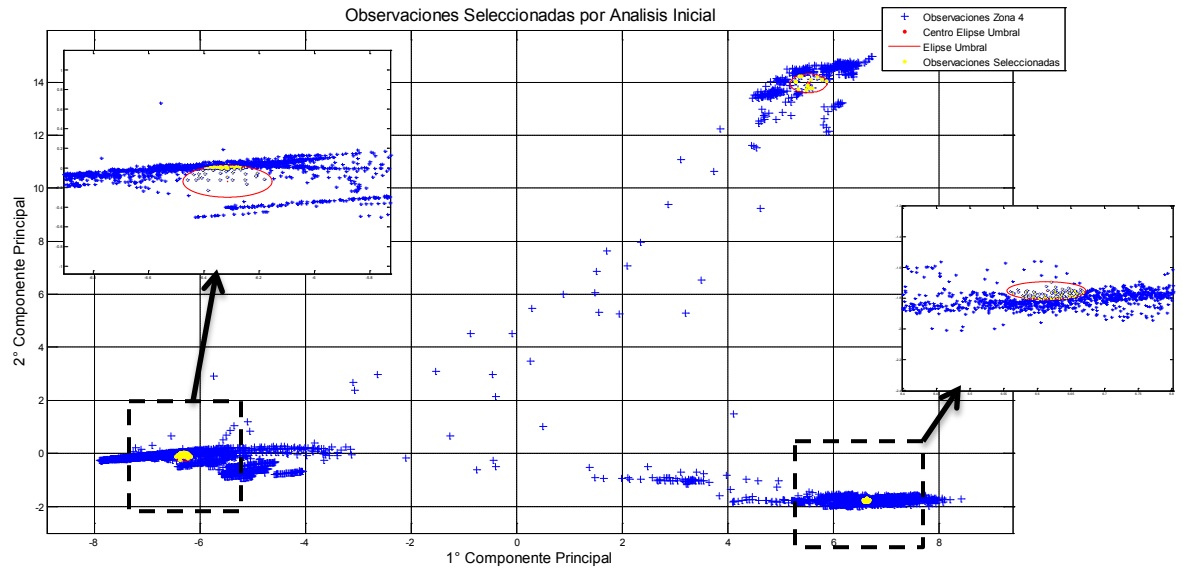


Figura 4.11: Observaciones seleccionadas por análisis inicial.

Luego de entregados los resultados obtenidos en la etapa de análisis inicial, se procede a entregar los gráficos correspondientes a cada una de las etapas posteriores. Básicamente se entregan gráficos que indican la calidad de las estimaciones realizadas por cada uno de los modelos SBM implementados a lo largo del algoritmo.

En primer lugar es necesario indicar que durante la evaluación de estos se determinó un comportamiento singular en los resultados que inducen a un permanente y localizado error en las estimaciones por el modelo SBM final. Este error es ilustrado en la Figura 4.12.

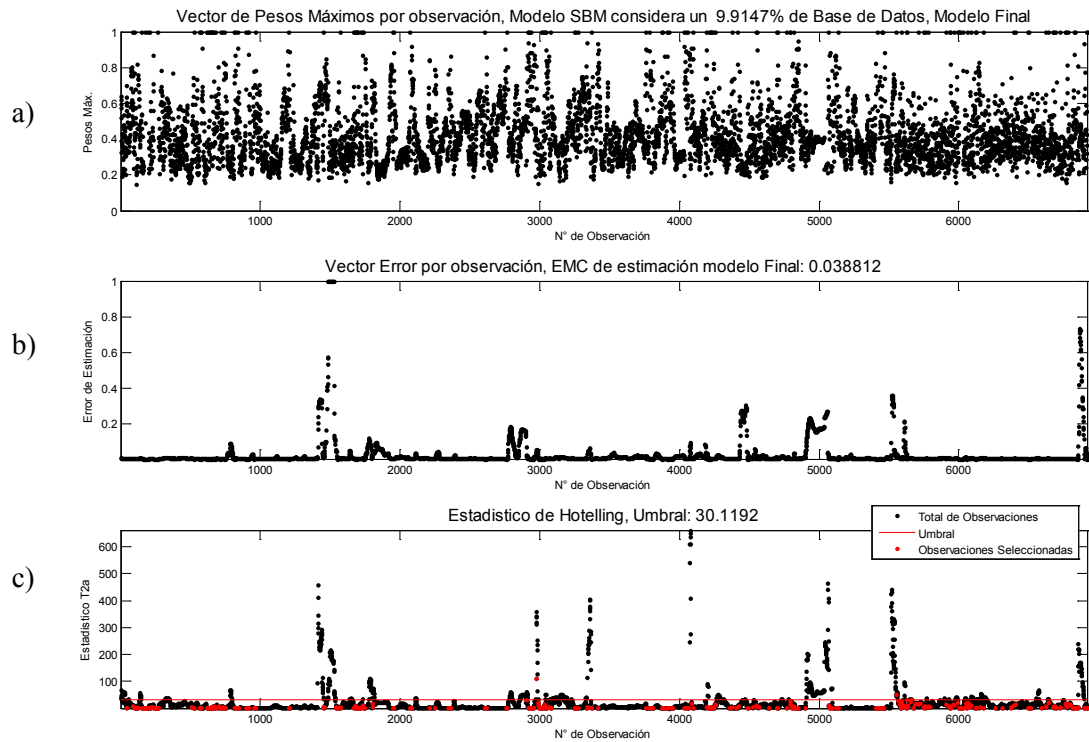


Figura 4.12: Resultado Preliminar obtenido con algoritmo de generación automática de modelos SBM.

Antes de describir tal comportamiento, se procederá a explicar la Figura 4.12, cuyo formato será utilizado reiteradamente para evaluar cada uno de los modelos implementados de aquí en adelante. Está constituido por 3 sub-gráficos: Vector de pesos máximos (Figura 4.12a), Vector Error por observación (Figura 4.12b) y Estadístico de Hotelling para las observaciones (Figura 4.12c). El primero de ellos tiene relación con la estructura del modelo SBM implementado, representa el vector de pesos máximos del valor absoluto de $\hat{\omega}$ (ver Ecuación (2.31)) y es un indicador de un reconocimiento exitoso de un vector de entrada en el conjunto de observaciones representativas que caracteriza un modelo SBM (C_{SBM}). Es así como, cuando para una observación el peso máximo sea 1 significa que tal entrada se encuentra en el conjunto C_{SBM} y, a medida que tal valor disminuye, indica que tal vector de entrada es cada vez más disímil a todos los elementos que constituyen C_{SBM} . La Figura 4.12b contiene el error de estimación del modelo SBM para cada una de las observaciones que constituyen la base de datos de entrenamiento, este error es calculado según Ecuación (3.3). Por último, la Figura 4.12c entrega el valor del estadístico de Hotelling para cada una de las observaciones, este estadístico es

calculado en base a la matriz de errores de estimación obtenidas por el modelo SBM; además se entrega el umbral de Hotelling (calculado según Ecuación (2.45)).

La irregularidad mencionada anteriormente puede ser visualizada en la Figura 4.12. Existe una zona en particular, localizada entre la observación 4900 y 5100 (ver Figura 4.13), que presenta un error de estimación permanente y comparativamente alto, siendo esto notorio tanto en el vector de pesos máximos, como en el vector de error de estimación y estadístico de Hotelling. Esta situación se presentó en cada uno de los modelos SBM implementados a lo largo del algoritmo de generación automática. La causa por la cual el algoritmo no incluyó ninguna de las observaciones de dicha zona se debe a que existe una restricción que impide agregar elementos que superen el umbral de Hotelling establecido por Ecuación (2.45). De esta manera, y dado que todos los elementos de tal zona se encuentran por sobre el umbral, el algoritmo negó la inclusión de algunos de las observaciones produciéndose un error permanente en la estimación de las salidas para dichos instantes.

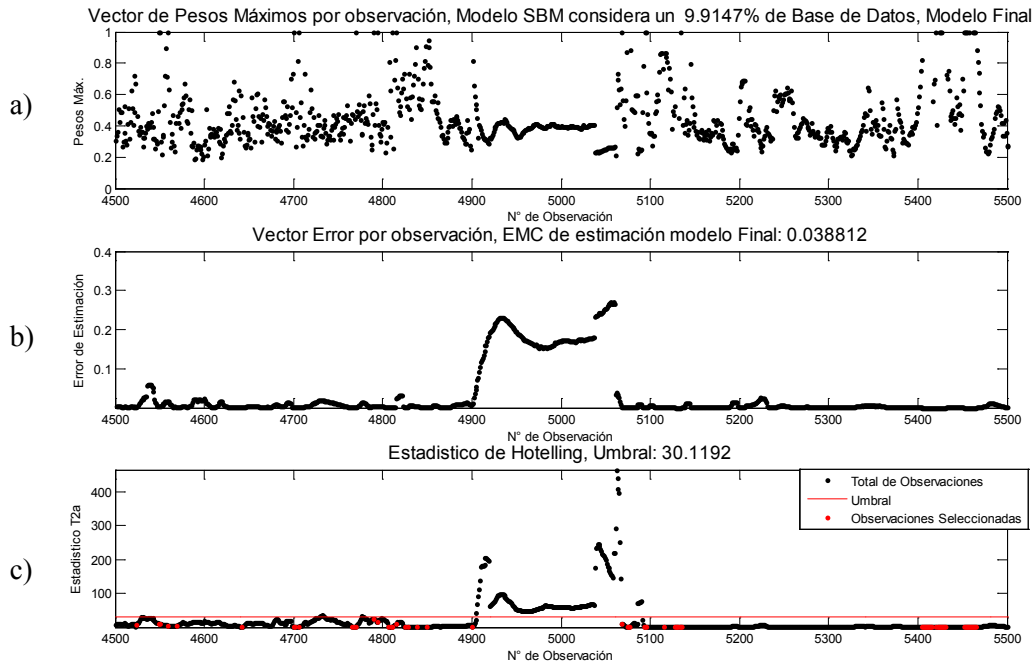


Figura 4.13: Aproximación a observación 5000 de la Figura 4.12.

El problema planteado puede ser solucionado al incluir una verificación en la etapa de selección de observaciones que se realiza vía el bloque K, descrito en sección 3.2.2.1.11. Esta verificación se basa en determinar el número de observaciones consecutivas que superan el umbral de Hotelling establecido. De esta manera, si se localizan 10 elementos consecutivos que superan tal umbral, serán agregados el quinto y décimo elemento de dicha serie siempre y cuando cumplan con los otros requerimientos establecidos por el algoritmo. Esta solución da flexibilidad a la inclusión de observaciones en el conjunto C_{SBM} y mejora las estimaciones del modelo SBM para tales zonas. Los resultados con la modificación introducida serán presentados en conjunto con la entrega de resultados finales.

A continuación se entrega una serie de gráficos que entrega el algoritmo de generación automática de modelos SBM y que describen la evolución de la calidad de los modelos SBM generados a lo largo de dicho proceso. Basado en lo explicado en Sección 3.2.1, se presentan una serie de figuras que representan cada una de las etapas del algoritmo.

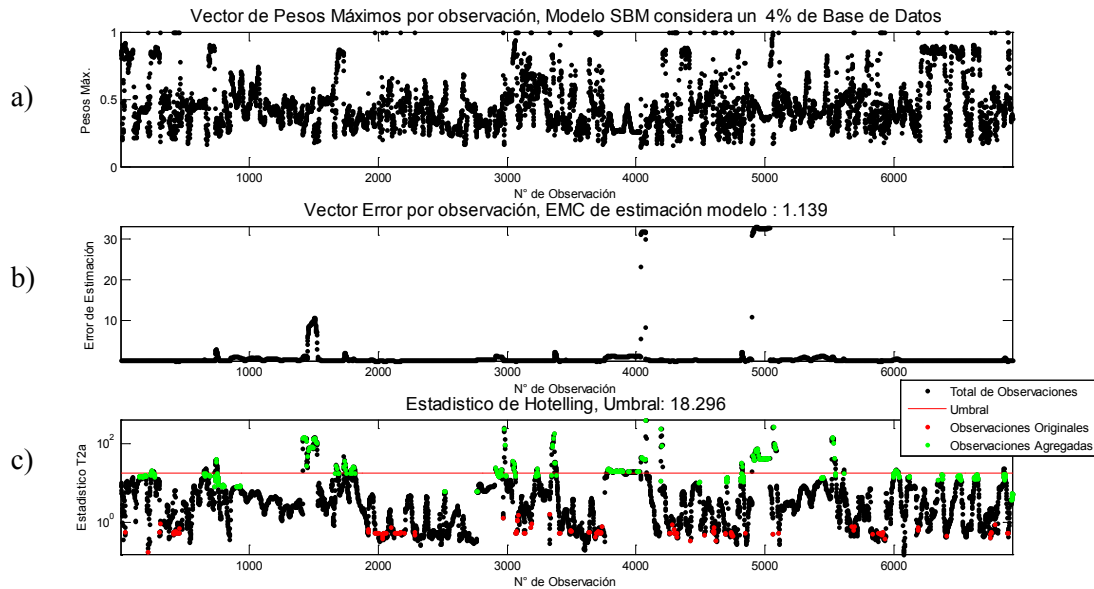


Figura 4.14: Resultados de modelo SBM inicial (4% de base de datos de entrenamiento). Gráfico c) en escala logarítmica.

La Figura 4.14 presenta los resultados obtenidos con el primer modelo SBM, el cual es creado en la etapa inicial del algoritmo. La Figura 4.14b se puede notar que existen errores de estimación puntuales y también existen zonas donde tal error es notablemente superior a las demás

estimaciones. Los errores de estimación puntuales no son considerados en los análisis debido a que pueden ser consecuencia de la dinámica no modelada del sistema o simplemente errores en los sensores. Por otro lado existen regiones que presentan EMC sostenidos, situación que puede deberse a dos motivos: i) las entradas del sistema no son representadas en el conjunto $C_{SBM_{inicial}}$, y por lo tanto el modelo SBM no es capaz de representar la salida, o ii) las entradas del sistema son similares a algunos elementos del conjunto $C_{SBM_{inicial}}$, pero el sistema actuó de manera diferente. Para determinar si el sistema se encuentra en alguno de tales escenarios la Figura 4.14a muestra los pesos máximos del valor absoluto de $\hat{\omega}$, el cual es un indicador de un reconocimiento del vector de entrada actual en el conjunto $C_{SBM_{inicial}}$. Tal como se mencionó anteriormente, un valor cercano a uno corresponde a observaciones que son similares a las consideradas en el conjunto $C_{SBM_{inicial}}$, en cambio, valores cercanos a cero indican que tales observaciones no se encuentran en la base del conjunto $C_{SBM_{inicial}}$ y por lo tanto presentan errores de estimación mayor. La Figura 4.14c provee información respecto al índice de Hotelling de cada observación, tal índice entrega una medida escalar del vector de errores de estimación obtenidos para cada una de las observaciones o instantes de tiempo. Este índice, y su correspondiente umbral de Hotelling, permite determinar cuáles son las observaciones que debiesen ser consideradas como elementos constituyentes de un posterior conjunto C_{SBM} que generará una versión mejorada del actual modelo SBM. Es así como; en rojo, se tienen las observaciones que constituyen el conjunto inicial $C_{SBM_{inicial}}$ y que a su vez construyen el modelo SBM inicial y; en verde, se presentan las observaciones que debiesen ser incorporadas a un posterior conjunto C_{SBM} . La Figura 4.14c permite visualizar el hecho de que han sido seleccionados (por el algoritmo) observaciones cuyo índice de Hotelling supera ampliamente el umbral establecido (línea roja), situación que se produce justamente debido a la modificación del algoritmo de manera de evitar el problema ilustrado en la Figura 4.12 y la Figura 4.13.

La argumentación anterior es también válida para los resultados mostrados en los gráficos posteriores (Figura 4.15, Figura 4.16, Figura 4.17, Figura 4.18, Figura 4.19 y Figura 4.20). Sin embargo, la principal característica diferenciadora es la disminución significativa tanto del EMC de las observaciones como del error del modelo calculado según Ecuación (3.7), para cada uno de los modelos creados en las diferentes etapas del algoritmo. Esta situación permite concluir de forma preliminar, que la modelación del proceso en estudio, en base al algoritmo diseñado, se ajusta paulatinamente al comportamiento medido en planta.

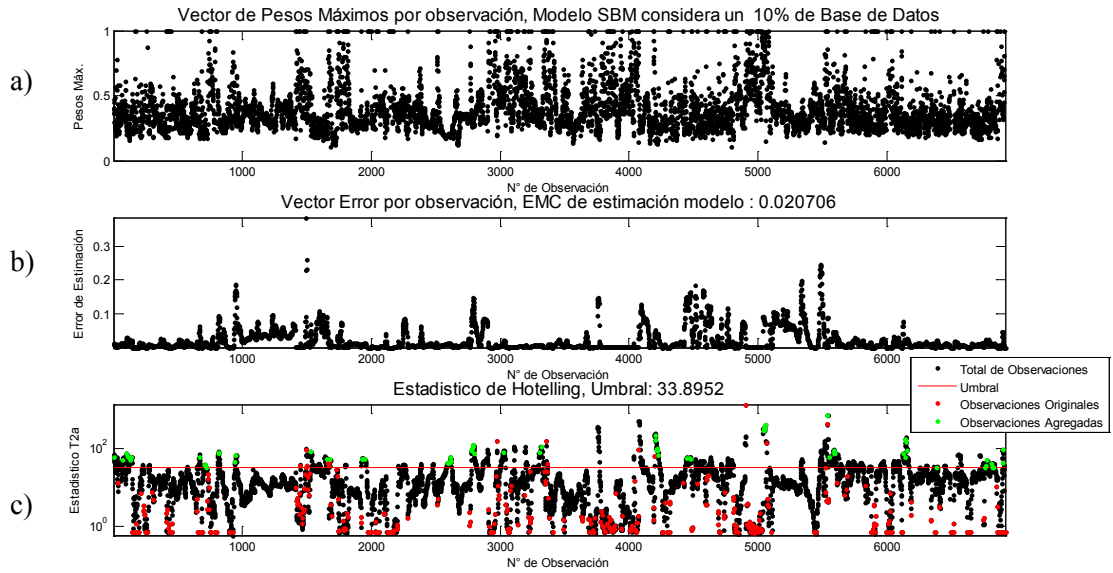


Figura 4.15: Resultados modelo SBM preliminar (10% de base de datos de entrenamiento). Gráfico c) en escala logarítmica.

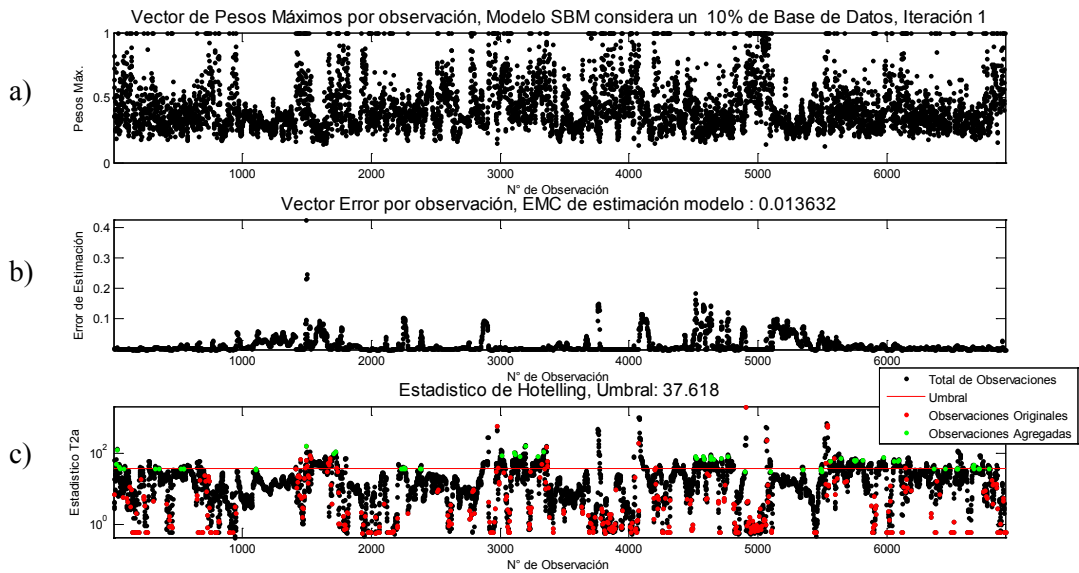


Figura 4.16: Resultado modelo SBM en 1º iteración. Gráfico c) en escala logarítmica.

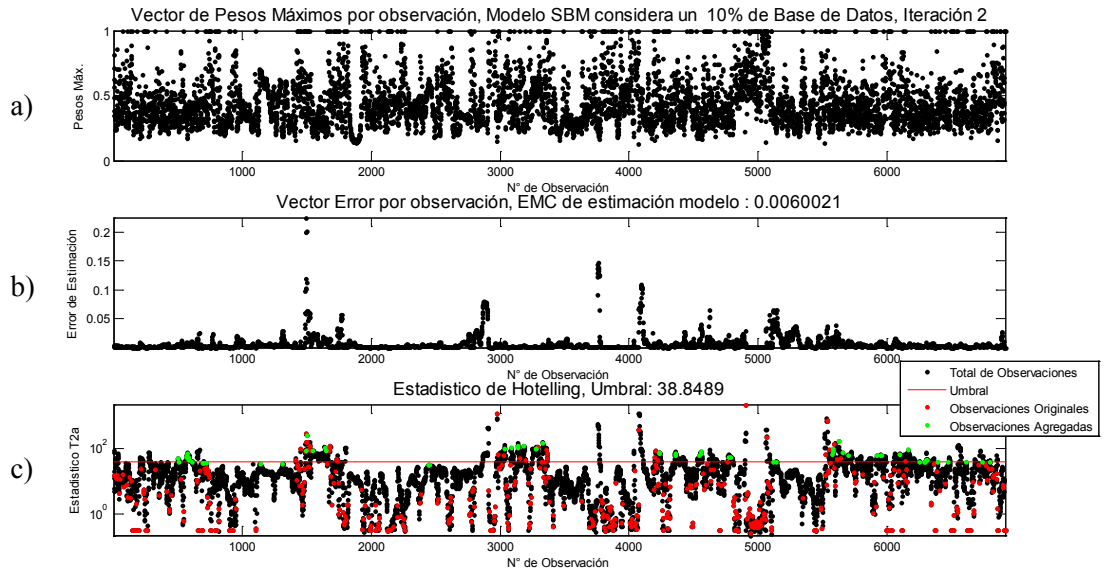


Figura 4.17: Resultado modelo SBM en 2^o iteración. Gráfico c) en escala logarítmica.

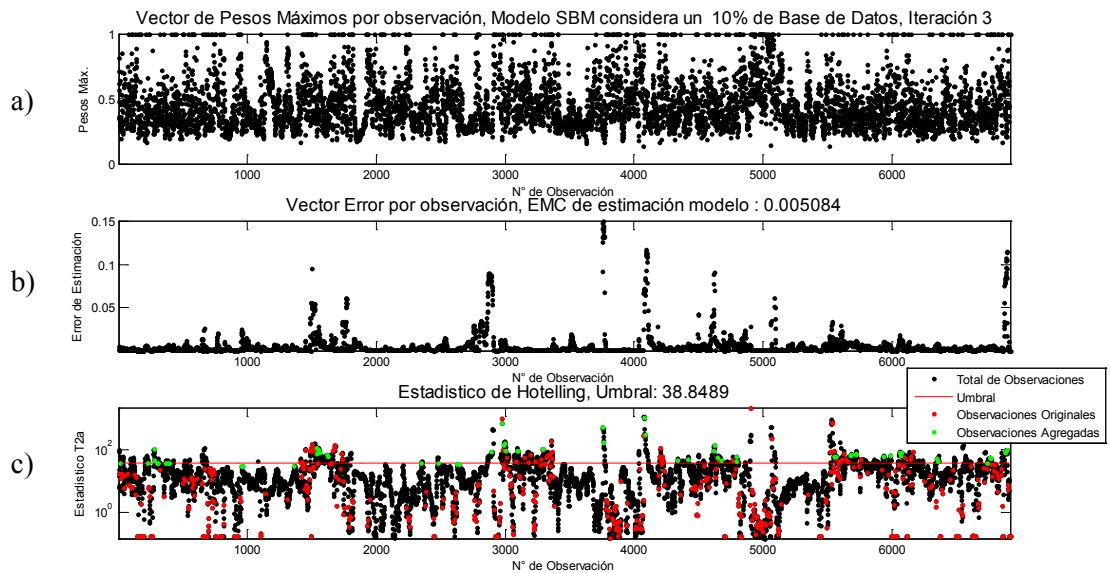


Figura 4.18: Resultado modelo SBM en 3^o iteración. Gráfico c) en escala logarítmica.

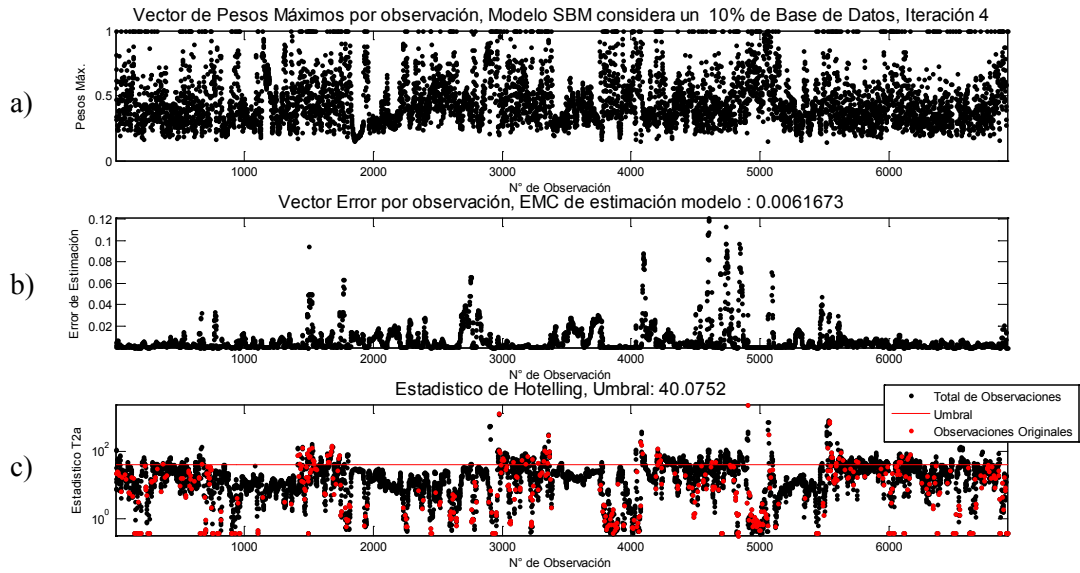


Figura 4.19: Resultado modelo SBM en 4^o iteración. Gráfico c) en escala logarítmica.

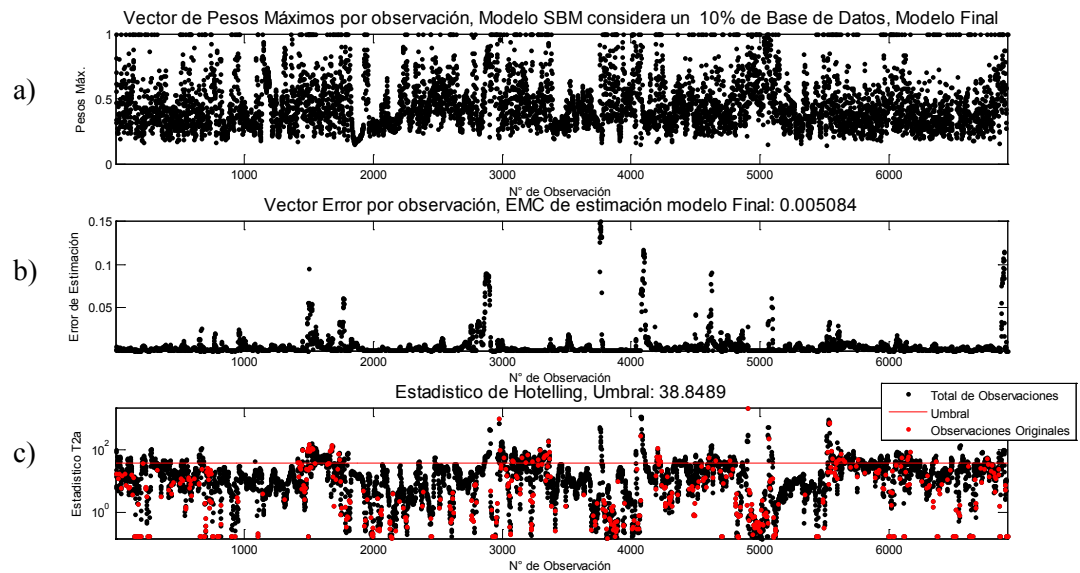


Figura 4.20: Resultado modelo SBM Final. Gráfico c) en escala logarítmica.

Luego de generado el modelo final, que representa al proceso según la base de datos de entrenamiento utilizada, existe una etapa opcional (Ajuste Final) en la cual se incorporan algunas observaciones extras de manera de disminuir o eliminar los valores más altos obtenidos en el EMC de las estimaciones y así, mejorar los resultados finales de la modelación. La Figura 4.21 corresponde a lo que el usuario vería al terminar el proceso iterativo de mejoramiento del modelo

SBM y en donde, al escoger la ejecución de la etapa de ajuste final, se procede a incorporar algunos instantes de tiempo que presentan EMC elevados y se procede a entregar una versión final del modelo SBM ajustado, cuyos resultados son los observados en la Figura 4.22.

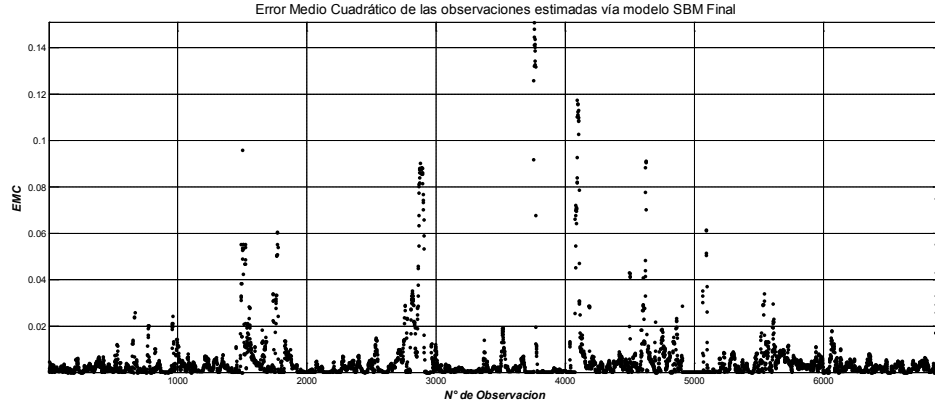


Figura 4.21: EMC (Error Medio Cuadrático) de las observaciones estimadas vía modelo SBM final.

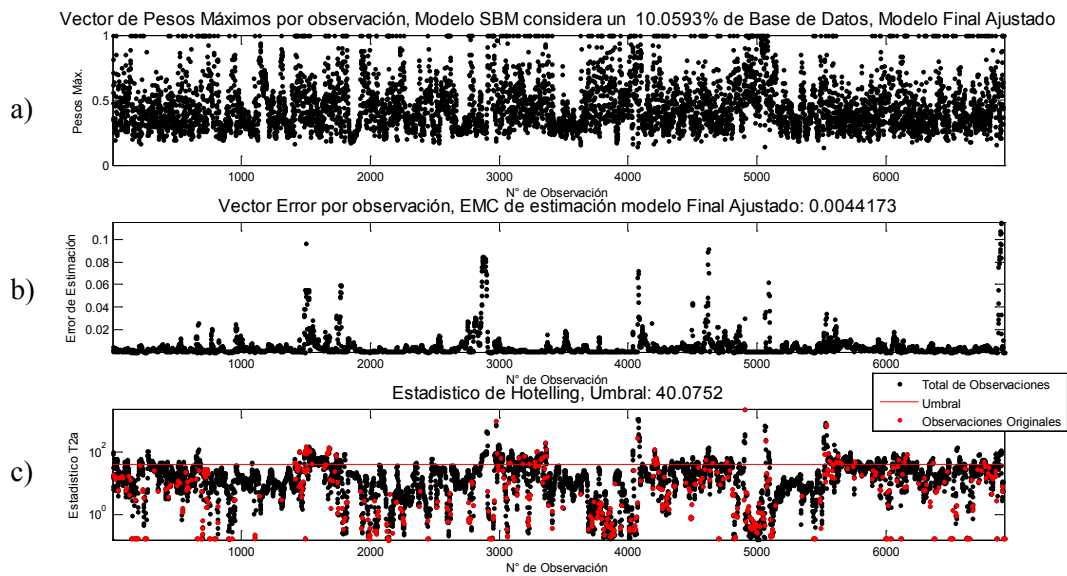


Figura 4.22: Resultado modelo SBM final ajustado. Gráfico c) en escala logarítmica.

La Figura 4.22 muestra que existe una disminución del error del modelo en comparación al modelo final, cuyos resultados se observan en la Figura 4.20. De esta forma, es posible ratificar un mejoramiento leve en la estimación de las salidas para la base de datos de entrenamiento utilizada.

La siguiente tabla entrega un resumen de los resultados obtenidos y en donde se puede observar la disminución del error del modelo para cada una de las etapas del algoritmo de generación automática, obteniendo un error de estimación del modelo final cercano al 0.5% del error obtenido en el modelo inicial. Además, tales datos están acompañados del tiempo acumulado de ejecución del algoritmo, lo cual permite corroborar que la generación de nuevos modelos de este tipo podrá ser posible en tiempo real en el caso de que se requiera, pues las bases de datos utilizadas están muestreadas a 10 [min].

	Modelo Inicial	Modelo Preliminar	Modelo Ite.1	Modelo Ite.2	Modelo Ite.3	Modelo Ite.4	Modelo Final	Modelo Final Ajustado
Error Modelo	1.1390	0.0207	0.0136	0.006	0.0051	0.0062	0.0051	0.0044
Tiempo de computo [s]	30.7	84.0	142.4	200.3	264.6	325.3	325.9	408.4

Tabla 4.3: Resumen de resultados obtenidos a lo largo de ejecución de algoritmo de generación automática de modelos SBM.

Es importante notar que el error obtenido por el modelo en iteración 4 fue superior al obtenido en la iteración 3. Esta situación ocurre debido a la manera en que se divide el conjunto inicial C1 descrito en la Sección 3.2.2.1.4, en la cual es realizado un ordenamiento descendente del vector estadístico de Hotelling obtenido con la primera matriz de error de estimación. Este ordenamiento produce que el último subconjunto que es eliminado contenga las mejores observaciones contenidas en C1 en términos de correcta estimación. Por lo tanto, no es extraño notar que al ser eliminados el error de estimación total del modelo empeore, en cuyos casos, se considera el modelo de la iteración 3 como el modelo final.

Entregados los resultados obtenidos en la generación automática de modelos SBM, es necesario utilizar algún criterio de comparación el cual permita ratificar que la modelación basada en similitud genera buenos resultados en contraste a otras técnicas de modelación. Es por este motivo que se diseña un modelo estático-lineal en los parámetros que permita realizar una diferenciación en base a los resultados de estimación de cada una de ellas. El criterio de comparación entre los modelos mencionados anteriormente estará basado en el error de modelación calculado según la Ecuación (3.7), siendo el mejor modelo aquel que tenga menor error.

Para la comparación entre modelos se definen 2 formas de dividir la base de datos T1 de manera de generar una base de entrenamiento y otra con fines de validación. En primer lugar, se realiza una división entre observaciones pares e impares, de esta manera se obtienen dos conjuntos con iguales condiciones de operación; sin embargo, dada la similitud entre muestras consecutivas, quizás los modelos aplicados a la base de datos de validación entreguen información redundante sobre su rendimiento. Por este motivo, la segunda división de T1 consiste en separarla de tal manera de construir el conjunto de entrenamiento con la primera mitad del total de observaciones, y por lo tanto, el conjunto de validación queda constituido con las observaciones restantes.

La modelación SBM está a cargo de los algoritmos propuestos en este trabajo y descritos en Sección 3.2. Mientras que la modelación multivariable lineal en los parámetros (de aquí en adelante, LMV) se realiza con la función de MATLAB® llamada *mvregress*, la cual asume una estructura de modelo indicada en Ecuación (2.3), Ecuación (2.10) y Ecuación (2.11). La estimación de parámetros es realizada a través del método de mínimos cuadrados descrito en la Sección 2.3.2.1. Las variables de entrada y salida son las indicadas en la Tabla 4.2.

A modo de introducir un factor que permita penalizar el número de parámetros m en relación al número N del conjunto de datos, se propone utilizar el índice *Akaike information Criterion* (AIC, [1]):

$$AIC = V \left(1 + \frac{2m}{N} \right) \quad (4.1)$$

Con V , una función de pérdida. En este caso, se utiliza la función indicada en Ecuación (3.7).

Esta penalización se realiza principalmente por la utilización de modelos lineales en los parámetros. Dado que los modelos SBM son clasificados como no-paramétricos, es posible considerar el número de observaciones del conjunto SBM (C_{SBM}) como parámetros estimados. De esta manera, dado que ambos tipos de modelos mejoran su estimación al aumentar el número de parámetros, se utiliza el criterio AIC y así, el modelo que obtenga un menor valor en el índice AIC será el que mejor represente el proceso en estudio.

Luego de lo anterior, se procede a entregar los resultados obtenidos tanto en entrenamiento como en validación para cada uno de los modelos:

T1: Pares (Entrenamiento)/Impares (Validación)

	SBM	LMV			
N° Variables Entrada	20	20			
N° Variables Salida	47	47			
N° Parámetros	346	987			
	Función de Perdida V: SBM	Función de Perdida V: LMV	N° observaciones	AIC-SBM	AIC-LMV
Entrenamiento	9.86E-03	4.11E-02	3460	1.18E-02	6.45E-02
Validación	1.01E-02	4.16E-02	3459	1.21E-02	6.54E-02

Tabla 4.4: Resultados de modelos SBM y LMV para base de datos T1 según división pares/impares.

Los resultados mostrados en la Tabla 4.4 permiten concluir que la técnica de modelación SBM presenta mejores estimaciones de las variables de salida que las realizadas por el modelo lineal en los parámetros, esto basado en que el error obtenido por parte del modelo SBM es aproximadamente 4 veces menor al error del modelo LMV, tanto en el conjunto de entrenamiento como en el de validación. Algo importante de notar es que los valores obtenidos tienden a ser similares en ambos conjuntos (entrenamiento y validación), debido a la forma en que estos fueron construidos, corroborando la necesidad de utilizar otra forma de separarlos para generar resultados completos.

La inclusión de un criterio de penalización por la cantidad de parámetros utilizados aumenta aún más la brecha entre ambos modelos, situación producida por el cálculo de más parámetros del modelo LMV con respecto al modelo SBM. La Figura 4.23 y la Figura 4.24 entregan los errores medios cuadráticos por observación, tanto en el conjunto de entrenamiento como en el conjunto de validación. Se puede visualizar que existen zonas en donde el error de estimación por parte del modelo lineal en los parámetros supera con creces al modelo SBM, debido principalmente a la inclusión de algunas observaciones en el conjunto C_{SBM} que producen error medio cuadrático igual a cero. Además, de los gráficos señalados se puede observar la similitud en la respuesta de ambos modelos frente a los conjuntos de entrenamiento y validación.

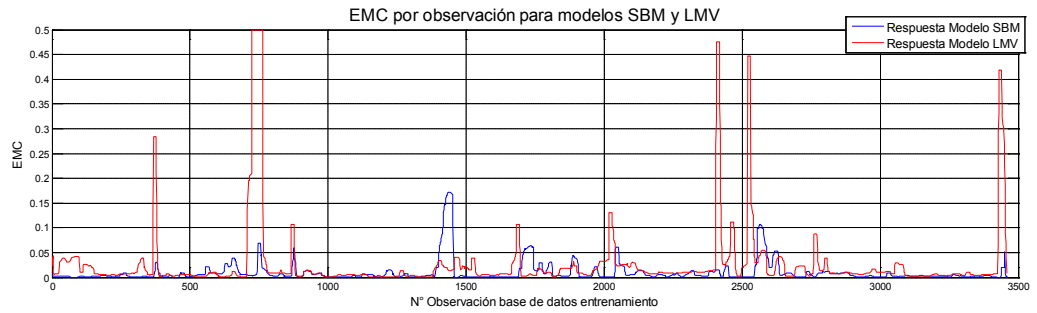


Figura 4.23: EMC por observación de modelos SBM y LMV para base de datos de entrenamiento 1. Datos saturados superiormente (máximo 0.5).

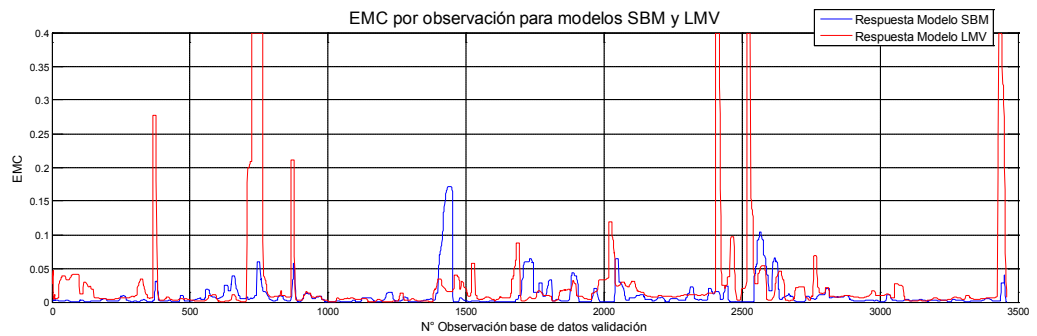


Figura 4.24: EMC por observación de modelos SBM y LMV para base de datos de validación 1. Datos saturados superiormente (máximo 0.4).

Dada la necesidad de comparar los modelos bajo todas las situaciones posibles, se generan los resultados de la modelación utilizando la segunda forma de separar la base de datos T1, los cuales se presentan a continuación:

T1: Observación 1-3460 (Entrenamiento) / 3461-6919(Validación)

	Modelo SBM	Modelo LMV			
N° Variables Entrada	20	20			
N° Variables Salida	47	47			
N° Parámetros	346	987			
	Función de Perdida V: SBM	Función de Perdida V: LMV	N° observaciones	AIC-SBM	AIC-LMV
Entrenamiento	5.74E-03	3.12E-02	3460	6.89E-03	4.90E-02
Validación	1.46E+00	2.22E+00	3459	1.75E+00	3.49E+00

Tabla 4.5: Resultados de modelos SBM y LMV para base de datos T1 según división 1-3460 y 3461-6919.

En esta oportunidad, existen condiciones de operación u observaciones en los datos de validación que no se asimilan a ninguno presente en la base de datos de entrenamiento. Esto debiese representar una gran dificultad para el modelo SBM, pues este se crea bajo el supuesto de que la información contenida en los datos de entrenamiento contiene todas las condiciones de operación del proceso. A pesar de lo mencionado, los resultados obtenidos permiten concluir que la modelación SBM genera mejores estimaciones de las salidas (con un error de estimación cercano a la mitad del obtenido con el modelo LMV) frente a observaciones no representadas en la base de datos de entrenamiento (ver Figura 4.25 y Figura 4.26), siendo esto importante si es utilizado en la herramienta de detección de anomalías, pues se obtiene menor índice de falsos positivos que los que se obtuviesen utilizando el modelo lineal en los parámetros.

Aun cuando los resultados permitan concluir que el modelo SBM representa más apropiadamente el proceso, se debe poner atención al error de estimación ya que aumenta alrededor de 2 órdenes de magnitud entre el conjunto de entrenamiento y el de validación, situación que empeora para las estimaciones cercanas a la observación 1500 de la Figura 4.26a. Esta situación deja en evidencia la importancia de utilizar una base de datos de entrenamiento rica en información de todas las condiciones de operación del proceso.

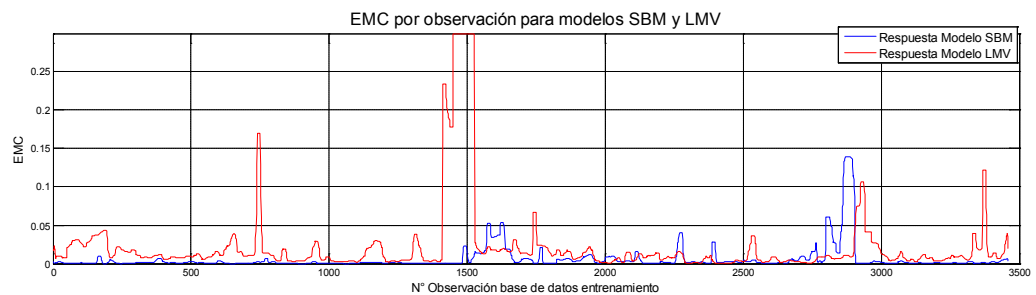


Figura 4.25: EMC por observación de modelos SBM y LMV para base de datos de entrenamiento 2. Datos saturados superiormente (máximo 0.3).

Finalmente, luego de entregados los resultados con respecto a la modelación SBM basada en lo descrito en la Sección 3.2, es posible concluir que ésta conduce a mejores estimaciones de las variables de salida en comparación a los modelos lineales en los parámetros. Siendo entonces, una técnica de modelación conveniente para ser utilizada en la herramienta de detección de anomalías propuesta en este trabajo y cuyos resultados se presentan en la próxima sección.

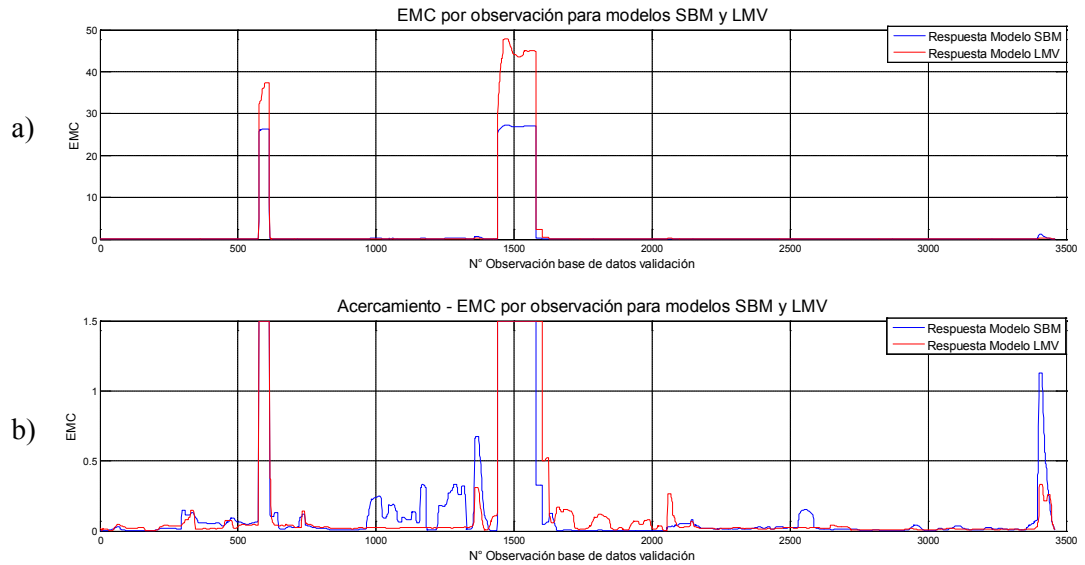


Figura 4.26: EMC por observación de modelos SBM y LMV para base de datos de validación 2 a) y acercamiento b). Datos Saturados superiormente (máximo 1.5).

4.3 RESULTADOS OBTENIDOS EN LA DETECCIÓN DE EVENTOS UTILIZANDO ALGORITMOS DE MODELAMIENTO BASADOS EN SIMILITUD Y BASE DE DATOS T2.

Luego de entregados los resultados sobre el algoritmo de modelación SBM, se procede a entregar una serie de resultados con respecto a la evaluación de la herramienta de detección de anomalías propuesta en la Sección 3.4.

La base de datos utilizada en esta etapa será la base T2 descrita en la Tabla 4.1. Esta base de datos tiene la particularidad de que se conocen los periodos de tiempo durante el cual el proceso operó en condiciones anómalas. A modo de entregar un análisis completo de los resultados generados se hace necesario describir las principales características de la anomalía en cuestión, basado en lo mencionado por operadores de planta e ingenieros especialistas. En primer lugar, el evento analizado corresponde a una acumulación de suciedad en zonas internas del compresor de la turbina, lo que se traduce en una disminución de la eficiencia de operación. Para solucionar este evento, se deben programar mantenciones extraordinarias que permitan limpiar tal zona,

debiéndose mantener apagada la turbina durante el lavado. Luego de eliminada tal acumulación, y según los operadores a cargo de la turbina, se produjo un aumento en la potencia máxima generada y en la presión del compresor, una disminución de la temperatura de descarga del compresor y una disminución del consumo de combustible (para una misma potencia).

La Figura 4.27 ilustra la potencia activa generada por el sistema turbina-generador. La anomalía denominada “suciedad en compresor” se presenta en dos oportunidades, una de ellas está acotada a las primeras observaciones, es decir, desde la observación 1 hasta 5400 aprox. (zona encerrada por elipse número 1 de color rojo). Luego de esto, se realiza una limpieza del compresor teniendo que detener la turbina, para luego operar en condiciones normales desde los instantes de tiempo 5650 hasta 11200 aprox. La segunda oportunidad en la que se presenta la anomalía en estudio se encuentra localizada entre las observaciones 11200 hasta 13200 aprox. (zona encerrada por elipse número 2 de color rojo). Posteriormente, se realiza una nueva limpieza entre los instantes 13600 y 16200, luego la turbina opera en condición normal hasta que se detiene el proceso en el instante 16200 aprox. Durante esta detención se realiza una mantención y renovación de algunos componentes de la turbina. Finalmente, se observa la operación de la turbina desde la observación 17500 hasta el último instante de tiempo registrado (zona encerrada por elipse número 3 de color gris).

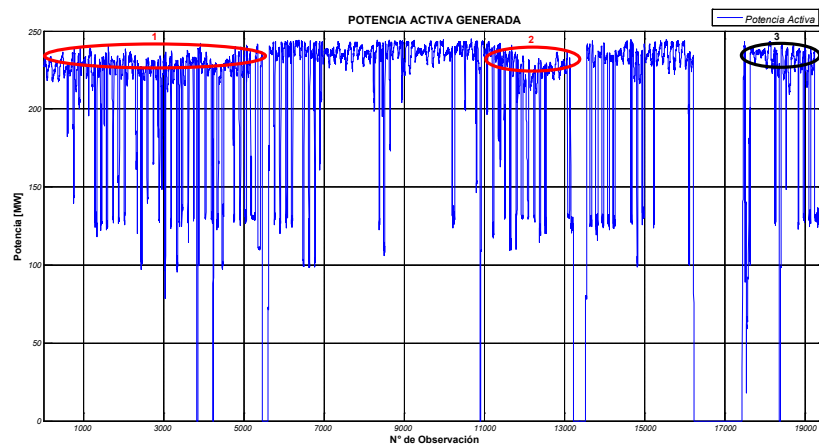


Figura 4.27: Potencia activa generada en sistema turbina-generador.

Para abordar el problema planteado se deberá en primer lugar, generar un modelo SBM utilizando como base de entrenamiento los instantes correspondientes a la 1ª zona de operación normal, es decir, entre las observaciones 5900 y 7489 exactamente. La metodología de estudio de eventos será utilizada para escoger las variables significativas que permiten detectar la anomalía en

cuestión. Para esta etapa se utilizará solamente las primeras 10000 observaciones. En otras palabras, se dejarán las observaciones 10001 en adelante para validar, tanto el modelo SBM para condición normal como la herramienta de detección de anomalías para condición anómala, dado que ambas condiciones se encuentran presentes en tal periodo de tiempo; además, la última parte de la base de datos T2 será utilizada para evaluar la robustez del modelo SBM debido a la modificación parcial del proceso.

	13	14	26	34	36	39	40	41
Variables de Entrada	42	58	61	63	64	76	86	101
N°	109	115	116	122	123	160	163	192
	193	194	195					
	4	5	6	7	8	9	10	11
	12	16	17	18	19	20	21	22
Variables de Salida	23	24	27	28	29	30	31	32
N°	33	7	38	46	55	68	69	70
	72	82	99	106	107	108	117	122
	123	126	127	128	176	177		

Tabla 4.6: Conjunto inicial de variables de entrada y salida.

La elección de las variables significativas empieza con una reducción de un conjunto de variables de entrada y salida inicial, que da como resultado lo indicado en la Tabla 4.6⁵. Con esto, se utilizó la metodología de estudio de eventos, que permitió reducir paulatinamente el número de variables de interés.

El procedimiento descrito en la metodología de estudio de eventos no sólo permitió encontrar las variables que producen una mejor modelación del proceso, y posterior detección de la anomalía. Adicionalmente, y luego de reiteradas iteraciones con diferentes niveles de reducción en la cantidad de variables incorporadas al modelo, se obtuvo un resultado particularmente llamativo, en donde se presentó un aumento considerable en el error de estimación del conjunto de salida (ver Figura 4.28, entre muestra número 2900 y 3700), situación por la cual fue necesario realizar un análisis en detalle de las variables involucradas de manera de localizar cual o cuales eran los factores que producían tal error. Cabe mencionar que, los instantes de tiempo en que el EMC se encuentra en el valor 50 se debe a que en tales instantes la turbina estuvo detenida, condición no

⁵ Las variables de entrada y salida utilizadas son identificadas sólo por su índice de la base de datos, de esta manera se respetan los contratos de confidencialidad del proyecto en el cual se enmarca este trabajo, detallado en Sección 1.4.

considerada en la base de datos de entrenamiento del modelo SBM, y que por lo tanto no necesita de mayor análisis.

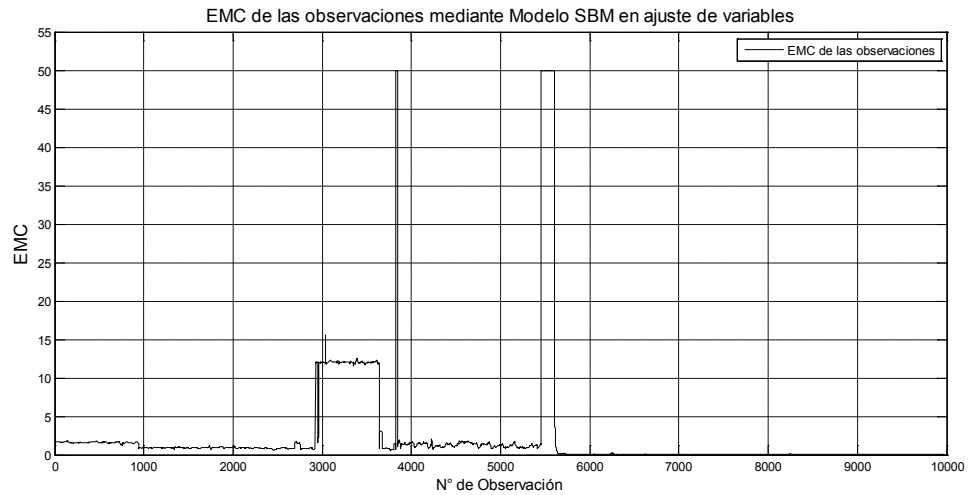


Figura 4.28: EMC de las observaciones mediante modelo SBM en proceso de ajuste de variables.

El resultado anterior indica que el modelo está estimando erróneamente todas las variables de salida o algunas de ellas son mal estimadas de manera significativa en la zona localizada entre las observaciones 2900 y 3700. Los estudios posteriores indicaron que la mala estimación del modelo se debió principalmente a tan sólo una variable de salida, la cual se ilustra a continuación:

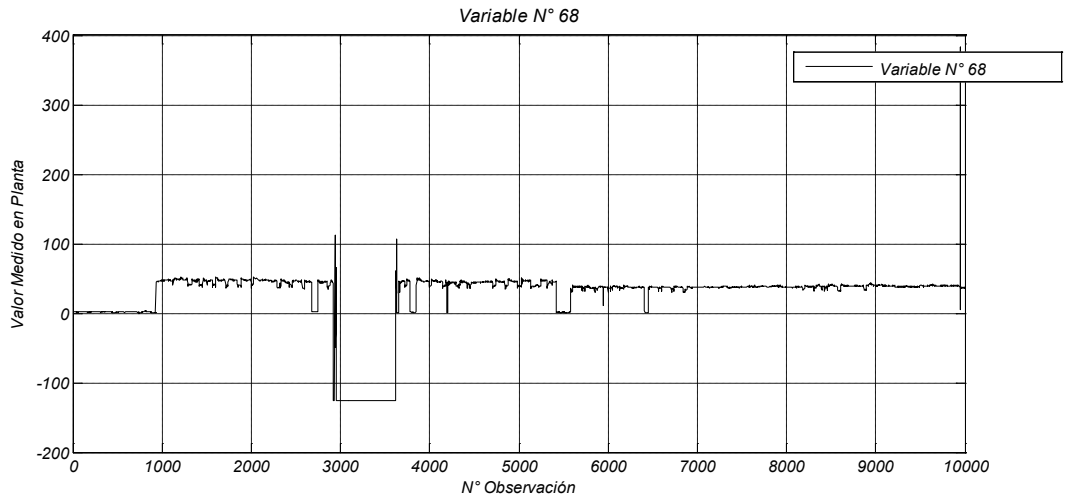


Figura 4.29: Variable N° 68 de la base de datos de entrenamiento.

Como se observa en la Figura 4.29, el modelo SBM obtenido a través de la metodología de estudio de eventos permitió identificar un comportamiento anormal de la variable 68,

correspondiente a la medición de un compuesto químico en la salida del proceso, donde se visualiza que entre los instantes 2900 y 3700 la señal parece descalibrarse y/o ajustarse a su valor más bajo. Esta situación fue descubierta a través del uso del algoritmo mencionado anteriormente y permite afirmar que el modelo SBM creado es, además, un buen indicador de las anomalías que puedan ocurrir en variables puntuales utilizadas por el algoritmo.

De la metodología de estudio de eventos y el algoritmo de generación automática de modelos, fue posible identificar las variables de mayor relevancia para la detección de la anomalía en análisis. Como se ha mencionado, el evento que se analiza se produce en la zona del compresor, es por ello que, las variables de interés y que se consideran en la creación del modelo final tienen que ver mayoritariamente con tal componente. Es importante recordar que, dada la naturaleza de los modelos basados en similitud, la idea central del sistema de detección es crear modelos que permitan representar condiciones de operación normal, así como también, generar modelos que posibiliten la detección e identificación de anomalías puntuales de un proceso industrial. Por lo tanto, el modelo final creado para la detección de este evento utiliza las siguientes variables para su creación:

Variable	TIPO
N° 58	ENTRADA - Flujo
N° 109	ENTRADA – Señal de Control
N° 115	ENTRADA – Temperatura 1
N° 122	ENTRADA – Posición Válvula 1
N° 123	ENTRADA – Posición Válvula 1
N° 38	SALIDA –Presión
N° 90	SALIDA – Potencia
N° 117	SALIDA – Temperatura 2

Tabla 4.7: Listado de variables de entrada y salida utilizadas en modelo SBM final.

El modelo SBM final que utiliza las variables indicadas en la Tabla 4.7 genera los siguientes resultados en la etapa de entrenamiento.

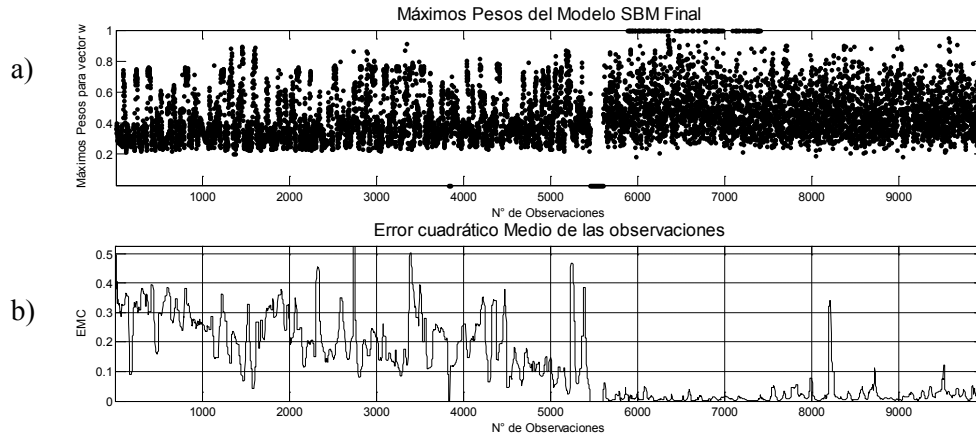


Figura 4.30: a) Vector de máximos pesos de ω y b) EMC de las observaciones según Modelo SBM Final.⁶⁻⁷

La Figura 4.30a entrega una visión de la estructura del modelo y de cómo estima las variables de salida, como se puede observar, la zona utilizada como base de entrenamiento del modelo (observaciones 5900 - 7489) contiene los más altos valores permitidos, debido a que son exactamente observaciones consideradas en el conjunto C_{SBM} . Es importante mencionar que tal conjunto corresponde al 10% del total de las observaciones de la base de datos de entrenamiento, es decir, considera aproximadamente 160 elementos con las cuales se estiman 10.000 observaciones.

Tanto el vector de pesos máximos del valor absoluto de $\hat{\omega}$ (Figura 4.30a), como el error cuadrático medio (Figura 4.30b), son indicadores complementarios que entregan información sobre la existencia de una anomalía en la zona localizada entre las observaciones 1 a 5500 aprox., zona etiquetada como operación anormal. Si bien, existen observaciones puntuales de la zona de operación normal que presentan un peso máximo bajo, el error medio cuadrático de tales instantes es comparativamente menor a la zona identificada como anómala. Esto se debe a que el vector de entrada para tales observaciones no es similar a algún vector del conjunto C_{SBM} , no obstante, una combinación lineal de todas las entradas de tal conjunto permite obtener una correcta estimación de las salidas. Situación que no ocurre para la zona anómala, pues teniendo observaciones con peso máximo bajo, el error medio cuadrático es comparativamente mayor a las zonas consideradas de operación normal. Además, de la Figura 4.30a es posible observar claramente los

⁶ En este gráfico y en los posteriores, la estimación del modelo para los instantes en que la máquina está detenida han sido llevadas a cero, de tal manera de no inferir en el análisis y/o discusiones.

⁷ Los gráficos obtenidos a través del modelo SBM han sido previamente filtrados para mejorar la lectura de los mismos.

instantes de tiempo en los cuales el proceso estuvo detenido, siendo esta condición de operación totalmente distinta a las observadas en el conjunto de entrenamiento, los pesos máximos obtenidos son cercanos o iguales a cero (observación cercanas a 4800 y 5500). Situación que en la Figura 4.30b no es notorio pues, en tales instantes, los valores EMC fueron llevados a cero para tener una mejor visualización de las observaciones en las zonas de interés (condición anormal y normal).

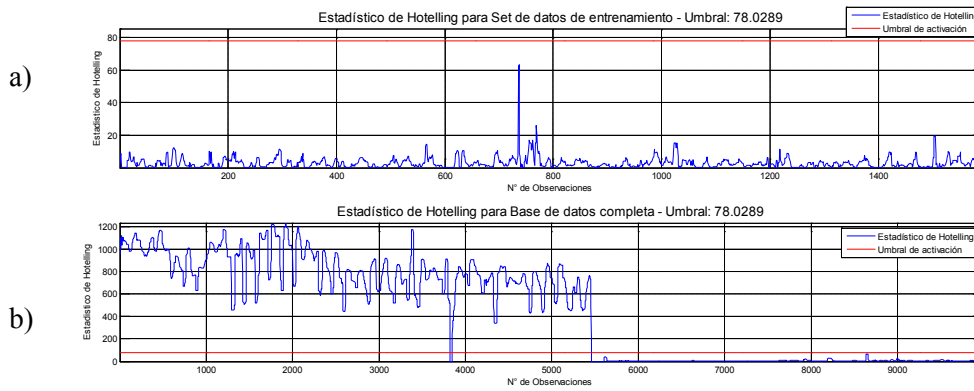


Figura 4.31: Resultados modelo SBM final, estadístico de Hotelling para las estimaciones.

La situación evidenciada en el gráfico anterior puede ser vista también en la Figura 4.31 donde se grafica el estadístico de Hotelling para el error de estimación, el cual es el indicador utilizado por la herramienta de detección de anomalías para activar la alarma de detección de la situación anómala.

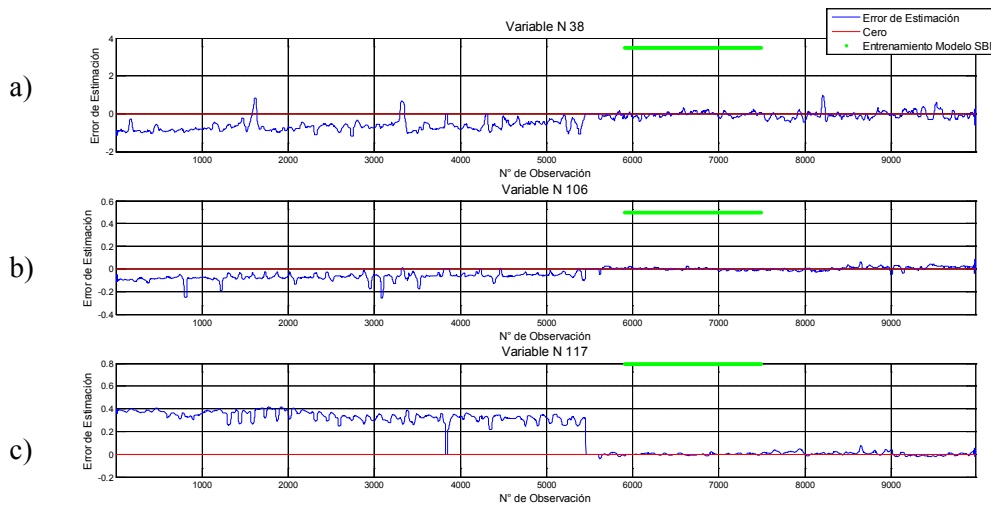


Figura 4.32: Error de estimación de las variables de salida del modelo SBM final para base de datos de entrenamiento.

La Figura 4.32 entrega el error de estimación de las variables de salida del modelo SBM creado, en donde se puede observar las dos zonas de interés, zona de anomalía (primeras 5500 observaciones aprox.) y zona normal. El desglose por variable de salida realizado permite al operador y/o usuario distinguir las variables conflictivas o que aportan más al error total del modelo. Esto puede ser utilizado para localizar la anomalía y así guiar la toma de decisión con respecto al tratamiento de la misma. La Figura 4.32a y Figura 4.32b tienen un comportamiento extraño en la zona comprendida entre las observaciones 5000 y 6000, esto es debido a que en aquella zona el proceso estuvo detenido y en tal caso, dado que el modelo SBM fue construido con una base de entrenamiento que contenía información del proceso sólo operando con gas, se producen errores de estimación mayores que no afectan a los resultados pues estas zonas desactivan la ejecución de la herramienta de detección.

Luego de aplicada la metodología de estudio de eventos a una base de datos que contenía información con respecto a la condición anómala y la condición normal, y de presentados los resultados acerca del mismo, se entrega a continuación los resultados obtenidos al evaluar la herramienta de detección de anomalías usando modelos SBM a la base de datos T2 completa. La idea principal es evaluar 3 situaciones claves:

- Validar el modelo SBM construido al analizar las observaciones de condición de operación normal después de la segunda limpieza del compresor.
- Evaluar la capacidad de detectar este tipo de anomalías al analizar la segunda aparición del evento según lo descrito en la Figura 4.27.
- Evaluar la robustez del modelo SBM implementado al analizar la operación del proceso con cambios en algunos componentes del mismo, detallado en la elipse número 3 de color negro de la Figura 4.27.

La Figura 4.33 ilustra el error de estimación del modelo SBM creado para cada una de las observaciones de T2.

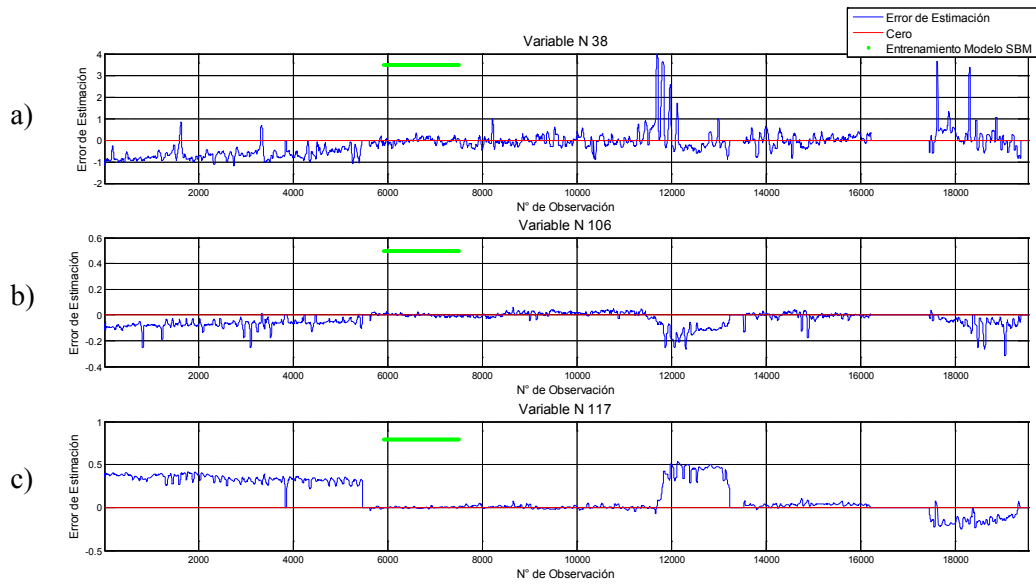


Figura 4.33: Error de estimación de variables de salida de modelo SBM final, base de datos T2.

En primer lugar, se puede observar que el modelo SBM creado permite estimar correctamente las variables de salida para la segunda zona de operación normal, notándose un comportamiento similar a lo obtenido para la primera zona de operación normal. Además, existe un aumento significativo en el error de estimación para las tres variables de salida en la segunda aparición del evento, esto es analizado a través del estadístico de Hotelling más adelante. Por último, se presenta un comportamiento extraño para los errores de estimación de la última zona de operación normal, esto se debe a las modificaciones hechas en el proceso y necesitan de un mayor análisis para una conclusión de la situación, análisis que será realizado al presentar la Figura 4.34: Indicador de anomalías según modelo SBM final, base de datos T2.

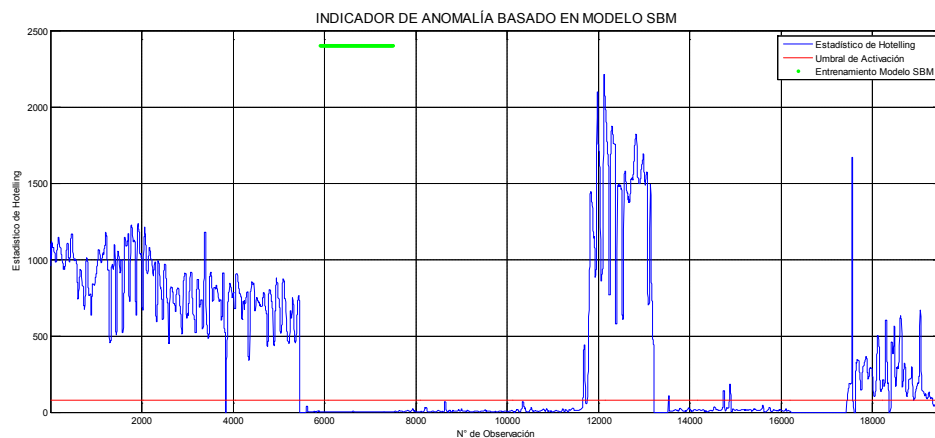


Figura 4.34: Indicador de anomalías según modelo SBM final, base de datos T2.

La Figura 4.34 corresponde al indicador de anomalías basado en el estadístico de Hotelling obtenido al analizar la matriz de errores de estimación. De la misma manera en que se explica la Figura 4.33, estos resultados permiten en primer lugar corroborar que existe una similitud en la respuesta del modelo para las primeras dos zonas de operación normal. Aun cuando existen observaciones, en la segunda zona de operación normal, cuyo índice de Hotelling supera el umbral de activación, estas son consideradas puntuales y pueden deberse a alguna medida errónea entregada por los sensores utilizados. Más aún, es posible establecer un criterio de activación de alarma que requiera de una sucesión de observaciones cuyo estadístico de Hotelling supere el umbral establecido.

Como puede observarse además, el índice de Hotelling en la zona en que ocurre el evento por segunda vez se eleva lo suficiente como para superar el umbral, siendo este comportamiento permanente a lo largo de la ocurrencia de dicho evento. Finalmente, se puede observar que para la tercera zona de operación normal, en donde se realizó un mantenimiento y renovación de elementos constituyentes de la turbina, el índice de Hotelling presenta valores lo suficientemente elevados y consecutivos que podrían provocar una activación errónea del evento. La causa principal del comportamiento obtenido radica en la naturaleza del modelo SBM, la cual hace uso de la información histórica del proceso para su construcción, creando mecanismos que relacionan las variables de entrada/salida y así, estimando las salidas a partir de nuevas medidas de entrada. En consecuencia, si se tiene un comportamiento del proceso diferente al encontrado en la base de datos de entrenamiento, la estimación y posterior detección de anomalías se verá afectada.

Para mejorar la estimación de las salidas y por ende, la detección exitosa del evento en estudio, se propone ajustar el modelo SBM creado de tal manera de incluir, en la base de datos de

entrenamiento, nuevas observaciones tomadas luego de la mantención/renovación de componentes de la turbina.

Siguiendo este esquema, se propone considerar dentro de la base de datos de entrenamiento, el periodo temporal localizado entre la observación 7800 y 8400 aprox., el cual es ilustrada en la Figura 4.35, donde se indica en color verde las zonas seleccionadas para ser consideradas en la generación del modelo SBM actualizado.

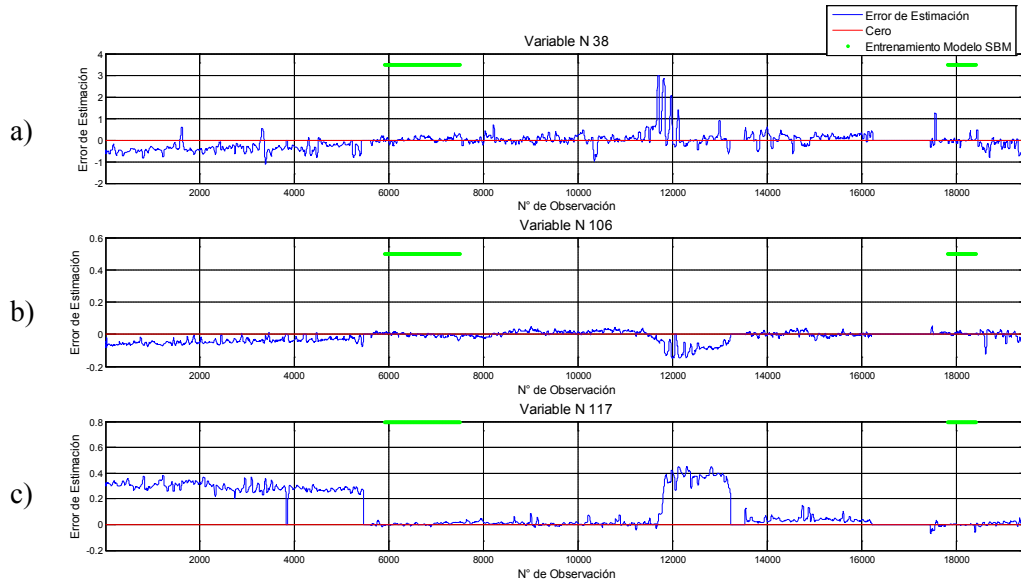


Figura 4.35: Error de estimación de variables de salida de modelo SBM actualizado, base de datos T2.

Las figuras anteriores entregan un detalle sobre los errores de estimación basados en el modelo SBM actualizado. En cada una de ellas se puede comprobar la mejoría en la estimación con respecto a lo obtenido con el modelo SBM Final (Figura 4.33), manteniendo las mismas características descritas para las zonas antecesoras. Es importante enfatizar que tanto en la estimación de la variable 106 y 117 (Figura 4.35a y Figura 4.35b, respectivamente) se presentan EMC máximos producto de la puesta en marcha de la turbina luego de la mantención/renovación de componentes constituyentes.

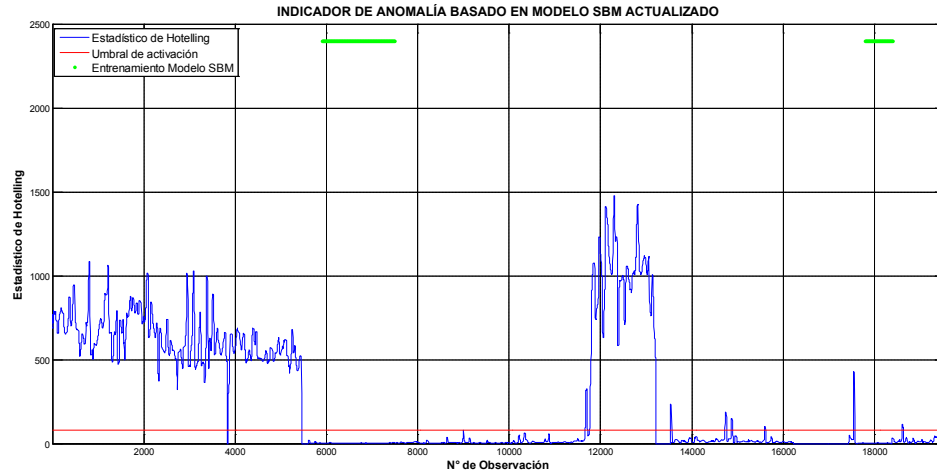


Figura 4.36: Indicador de anomalías según modelo SBM actualizado, base de datos T2.

La Figura 4.36 muestra los resultados utilizando el modelo SBM actualizado, en esta oportunidad se puede observar con mayor claridad la mejoría del indicador de Hotelling con respecto a los resultados en la Figura 4.34 manteniéndose el mismo comportamiento para las dos primeras zonas de condición normal y para las 2 zonas de operación anormal. Además, es posible afirmar que para los últimos instantes de la base de datos, el estadístico de Hotelling se ubica bajo el umbral de activación definido. Aun cuando se identifican puntos particulares en los cuales se sobrepasa el umbral, que son provocados por el comportamiento variante de la potencia activa de la turbina (producto de la puesta en marcha de la turbina luego de las modificaciones realizadas) y por situaciones anormales en la variable 38 y 117 no registradas con anterioridad, es posible concluir que el modelo SBM actualizado vuelve a ser un buen indicador de anomalías en el compresor.

Para finalizar la entrega de resultados y, basados en la necesidad de elaborar una herramienta de detección temprana de eventos anómalos ubicados en la zona del compresor, se presentan dos gráficos (Figura 4.37 y Figura 4.38) en los cuales se ilustran tanto la potencia activa de la máquina como el indicador de Hotelling resultante basado en el modelo SBM actualizado.

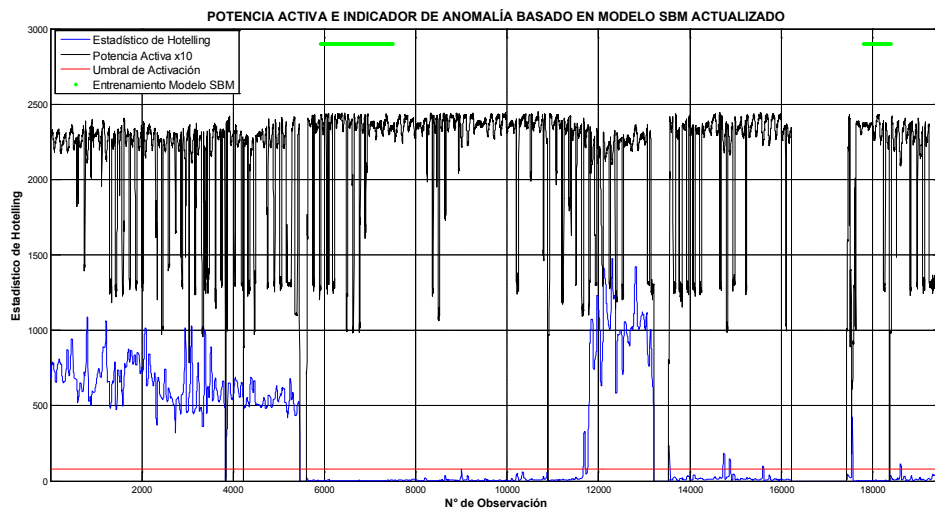


Figura 4.37: Comparación entre potencia generada e indicador de anomalías para modelo SBM actualizado.

Dado el nivel de dificultad para leer el gráfico anterior, se hace necesario entregar una vista más detallada del mismo. En particular, interesa conocer el momento exacto en el cual ocurre el evento suciedad de compresor y el instante en donde el indicador de anomalías diseñado detecta la situación. Por esto, se hace un acercamiento a la zona comprendida entre la observación 11500 y 12200.

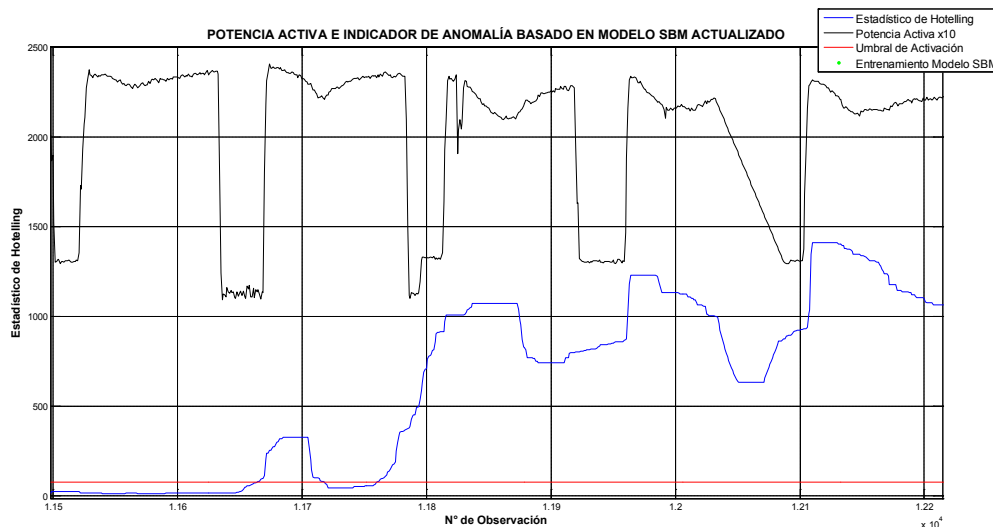


Figura 4.38: Acercamiento a zona de interés de Figura 4.37.

Este último gráfico permite conocer de forma más exacta el momento en el cual el estadístico de Hotelling, obtenido del modelo SBM actualizado, supera el umbral de activación propuesto. Para

obtener una conclusión concreta sobre la detección temprana del evento para este proceso en particular, es posible contrastar esta información con aquella obtenida por parte de los operadores y/o especialistas del mismo. Sin embargo, dado que no existe una herramienta de detección implementada en tal proceso, este primer enfoque entrega un indicador válido, en línea y comprobado sobre la condición de operación del proceso.

CAPÍTULO 5. CONCLUSIONES

La detección de anomalías es un problema importante que ha sido explorado en diversos campos de investigación y áreas de aplicación. El interés principal en el desarrollo de estas herramientas está relacionado con la reducción de costos —tanto humanos como económicos— en la industria contemporánea. En base a esto, se propuso el desarrollo de una herramienta de detección de anomalías para procesos industriales multivariados usando una técnica de modelamiento basada en similitud. Para lograr tal propósito, se presentaron una serie de objetivos específicos que permitieron, en su conjunto, cumplir con el objetivo planteado. Sobre estos objetivos se puede mencionar lo siguiente:

Usando el enfoque de detección de anomalías basada en observadores, se diseñó un algoritmo de generación automática de modelos basados en similitud. La idea general era disponer de una herramienta iterativa que permita generar este tipo de modelos en base a un conjunto de datos de entrenamiento. Los resultados obtenidos al evaluar la capacidad de modelamiento fueron analizados a lo largo de la evolución del algoritmo iterativo, y luego comparados con otro tipo de modelación, exactamente con los modelos lineales en los parámetros. Los primeros análisis permiten señalar que la evolución del algoritmo, y de los modelos SBM propuestos, mejoraron paulatinamente a lo largo de su ejecución, reduciendo el error de estimación a cerca del 0.5 % del error de estimación del modelo SBM inicial. Las comparaciones realizadas a la capacidad de modelamiento de la herramienta entregaron resultados que permiten concluir que el algoritmo diseñado genera modelos que representan mejor al proceso. Esta conclusión tiene su respaldo en el bajo error de estimación obtenido por el modelo SBM, del orden del 25% del obtenido con los modelos lineales en los parámetros. Cabe mencionar que los modelos tratados representan el proceso desde el punto de vista estático, por ello queda como trabajo propuesto desarrollar un mecanismo que permita generar modelos SBM dinámicos. Para emular el comportamiento dinámico de un sistema podrían ser incluidas mediciones pasadas (tanto entradas como salidas) como regresores para predecir la salida.

Los resultados de la etapa de prueba de la herramienta de detección de anomalías usando modelos basados en similitud, permiten corroborar la utilidad de contar con una metodología de estudio de evento. Esta metodología permitió detectar correctamente las variables que no aportan al objetivo de detección, es más, se utilizó principalmente para eliminar variables que conducen a resultados

poco efectivos en lo que a detección de anomalías se refiere, causado por la nula relevancia de tales variables en el comportamiento anormal observado. Además, se logró identificar las variables más significativas para el proceso de detección del evento estudiado.

Los resultados conseguidos en la detección de un evento en particular, usando la herramienta de detección de anomalías, fueron analizados desde tres enfoques distintos los cuales se detallan a continuación:

En primer lugar, se validó el modelo SBM generado por la herramienta de detección de anomalías. Esta validación fue lograda al corroborar la correcta estimación de las variables de salida del modelo frente a observaciones del proceso no consideradas en la base de datos de entrenamiento. Además, esta situación tiene el respaldo del indicador de anomalías utilizado, el cual estuvo siempre por debajo del umbral de activación del evento para tales zonas.

En segundo lugar, el modelo SBM construido y utilizado en la herramienta de detección de anomalías resultó ser eficaz en la detección, pues en las zonas caracterizadas por la ocurrencia del evento anómalo el indicador de anomalías superó el umbral de activación del mismo.

Como tercer punto importante, se analizaron los resultados obtenidos luego de una intervención mayor del proceso; el cual consistió en una renovación de elementos constituyentes del mismo. El objetivo de este análisis fue evaluar la robustez del modelo SBM creado para la estimación. Los resultados obtenidos dieron evidencia de una situación no deseada, en donde para la operación normal del proceso, el indicador de anomalía presentó valores lo suficientemente elevados y consecutivos que provocan una activación errónea del evento. Esta situación evidenció el impacto en la eficacia de la detección de anomalías al intervenir el proceso. No obstante, un reajuste el modelo SBM, el cual consideró utilizar observaciones del proceso luego de la intervención del mismo, produjo una mejoría en la herramienta de detección diseñada, lográndose comportamientos deseados y vistos en las anteriores zonas de la base de datos.

Con respecto a la herramienta de detección de anomalías, queda por trabajo a futuro la implementación de una interfaz que permita utilizar tal herramienta dentro del sistema de control y monitoreo de un proceso industrial. Además, considerando un eventual desarrollo del algoritmo de generación automática de modelos SBM dinámicos, es lógico pensar en la utilización de tal actualización en el análisis de anomalías que necesitan de un mecanismo dinámico para su detección.

BIBLIOGRAFÍA

1. Akaike, H., "A new look at the statistical model identification," *Automatic Control, IEEE Transactions on*, vol.19, no.6, pp. 716- 723, Dec 1974.
2. Araya. J. (1997). Estimación de parámetros para su aplicación a diagnóstico de fallas. (Tesis de Magister – Pontificia Universidad Católica de Chile), [en línea]. Disponible en: <http://www2.ing.puc.cl/~jfaraya/thesis/Tesis.pdf> [Consulta: 13 octubre 2011].
3. Baskiotis, C.; Raymond, J.; Rault, A. , "Parameter identification and discriminant analysis for jet engine mechanical state diagnosis," *Decision and Control including the Symposium on Adaptive Processes, 1979 18th IEEE Conference on*, vol.18, pp.648-650, Dec. 1979.
4. Beale, G.O.; Kim, J.H., "Fisher discriminant analysis and the T2 statistic for process fault detection and classification," *IECON 02 [Industrial Electronics Society, IEEE 2002 28th Annual Conference]*, vol.3, pp. 1995- 2000, 5-8 Nov. 2002.
5. Chandola V., Banerjee A., Kumar V., "Anomaly Detection: A Survey", *ACM Computing Surveys*, Vol. 41(3), Article 15, July 2009.
6. Chen, A.; Elsayed, E.A., "Mean estimate for Shewhart-chart-monitored processes subject to random shifts," *Systems, Man, and Cybernetics*, 1998. 1998 IEEE International Conference on, vol.5, pp.4687-4692, 11-14 Oct 1998.
7. Chiang, L.H.; Russell E.L. and Braatz R.D., *Fault detection and diagnosis in industrial systems*, Great Britain: Springer, 2001, pp. 21-24.
8. Chiang, L.H.; Russell E.L. and Braatz R.D., *Fault detection and diagnosis in industrial systems*, Great Britain: Springer, 2001, pp. 35-38.
9. Chiang, L.H.; Russell E.L. and Braatz R.D., *Fault detection and diagnosis in industrial systems*, Great Britain: Springer, 2001, pp. 71-74.
10. Chiang, L.H.; Russell E.L. and Braatz R.D., *Fault detection and diagnosis in industrial systems*, Great Britain: Springer, 2001, pp.3-9.
11. Clark, R.N.; Fosth, D.C.; Walton, V.M., "Detecting Instrument Malfunctions in Control Systems," *Aerospace and Electronic Systems, IEEE Transactions on*, vol.AES-11, no.4, pp.465-473, July 1975.
12. Deckert, J.; Desai, M.; Deyst, J.; Willsky, A. , "F-8 DFBW sensor failure identification using analytic redundancy," *Automatic Control, IEEE Transactions on*, vol.22, no.5, pp. 795- 803, Oct 1977.
13. Doyle, R.; Charest, L. Jr.; Rouquette, N.; Wyatt, J., "Causal Modeling and Event-driven Simulation for Monitoring of Continuous Systems," *Jet Propulsion Laboratory - Technical Report Server*, Oct. 1993.
14. Duda, R.; Hart, P.; Stork, D., *Pattern Classification*, New York: John Wiley & Sons, 1999, pp. 114-121.
15. El-Shal, S.M.; Morris, A.S., "A fuzzy expert system for fault detection in statistical process control of industrial processes," *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, vol.30, no.2, pp.281-289, May 2000.

16. Fan, M.Jr.; Guo, R.S.; Chang, S.H.; Lee, J.H., "Abnormal trend detection of sequence-disordered data using EWMA method [wafer fabrication]," Advanced Semiconductor Manufacturing Conference and Workshop, 1996. ASMC 96 Proceedings. IEEE/SEMI 1996, pp.169-174, 12-14 Nov 1996.
17. Frank, P.M., "Fault diagnosis in dynamic systems using analytical and knowledge-based redundancy - a survey and some new results," *Automatica*, vol.26, no.3, pp. 459-474, May 1990.
18. Fuente, M.J.; Garcia-Alvarez, D.; Sainz-Palmero, G.I.; Villegas, T. , "Fault detection and identification method based on multivariate statistical techniques," Emerging Technologies & Factory Automation, 2009. ETFA 2009. IEEE Conference on, pp.1-6, 22-25 Sept. 2009.
19. Gazzana, D.S.; Oliveira, M.O.; Bretas, A.S.; Lerm, A.A.P.; Bettiol, A.L.; Da S Goncalves, M.A. , "An expert system for substation fault detection in thermoelectric generation plants," Modern Electric Power Systems (MEPS), 2010 Proceedings of the International Symposium , pp.1-6, 20-22 Sept. 2010.
20. Gong, L.; Schonfeld, D., "Space Kernel Analysis," Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on, pp.1577-1580, 19-24 April 2009.
21. González, G.D.; Paut, R.; Cipriano, A.; Miranda, D.R.; Ceballos, G.E. , "Fault detection and isolation using concatenated wavelet transform variances and discriminant analysis," Signal Processing, IEEE Transactions on , vol.54, no.5, pp. 1727- 1736, May 2006.
22. Hao Long; Xinmin Wang, "Aircraft fuel system diagnostic fault detection through expert system," Intelligent Control and Automation, 2008. WCICA 2008. 7th World Congress on, pp.7104-7107, 25-27 June 2008.
23. Huang, S.; Tan, K.K., "Fault Detection and Diagnosis Based on Modeling and Estimation Methods," Neural Networks, IEEE Transactions on, vol.20, no.5, pp.872-881, May 2009.
24. Isermann R., *Fault-Diagnosis Systems: An Introduction from Fault Detection to Fault Tolerance*. Heidelberg: Springer, 2006, pp. 20-21.
25. Isermann R.; Ballé P., "Trends in the Application of Model-Based Fault Detection and Diagnosis of Technical Processes". *Control Engineering Practice*, vol. 5, no. 5, pp. 709-719, May 1997.
26. Isermann, R., "Process fault detection based on modeling and estimation methods – A survey," *Automatica*, vol.20, no.4, pp.387-404, July 1984.
27. Kabbaj, N.; Ramzi, M.; Dahhou, B.; Youlal, H.; Enea, G., "Fault detection and isolation in a greenhouse using parity relations," Emerging Technologies and Factory Automation, 2003. Proceedings. ETFA '03. IEEE Conference, vol.2, pp. 747- 752, 16-19 Sept. 2003.
28. Liang, Z.; Cao J.; Zhou J. , "A statistical method for health diagnosis of concrete bridge based on EWMA control chart and reliability analysis," Fuzzy Systems and Knowledge Discovery (FSKD), 2010 Seventh International Conference on , vol.6, pp.2584-2588, 10-12 Aug. 2010.
29. Ljung, L., *Modeling of Dynamic System*, New Jersey: Prentice Hall, 1994, pp. 19-21.
30. Ljung, L., *System Identification, Theory for the user*, New Jersey: Prentice Hall, 1999, caps 1-4.
31. Moody, J.; Darken, C., "Learning with localized receptive fields," Proceedings of the 1988 Connectionist Models Summer School, eds. Touretzky, Hinton, and Sejnowski. Morgan-Kaufmann, Publishers, 1988.
32. Moran, M.J., *Engineering Thermodynamics Mechanical Engineering Handbook*, Boca Raton: CRC Press LLC, 1999.

33. Moreno L.; Garrido S.; Balaguer C., Ingeniería de control. Modelado y control de sistemas dinámicos. Barcelona: Ariel, 2003, pp. 238-245.
34. Oppenheim, A. V.; Willsky, A. S., Signal and Systems, New Jersey: Prentice Hall, 1996, pp. 38-56.
35. Patton, R.J.; Chen, J., "Robustness in quantitative model-based fault diagnosis," Intelligent Fault Diagnosis - Part 2: Model-Based Techniques, IEE Colloquium on, vol.4, no., pp. 1-417, 26 Feb 1992.
36. Schimek, M., Smoothing and regression: Approaches, computation, and application, New York: John Wiley & Sons, 2000.
37. Senouci, K.; Bendaoud, A.; Medles, K.; Tilmatine, A.; Dascalescu, L. , "Comparative Study between the Shewhart and CUSUM Charts for the Statistic Control of Electrostatic Separation Processes," Industry Applications Society Annual Meeting, 2008. IAS '08. IEEE , pp.1-5, 5-9 Oct. 2008.
38. Shewart W.; Wilks S., Methods and applications of linear models. New Jersey: John Wiley & Sons, 2003, pp. 1-20.
39. Sjöberg, J.; Zhang Q.; Ljung L.; Benveniste A.; Delyon B.; Glorennec P.; Hjalmarsson H.; and Juditsky A. , "Nonlinear black-box modeling in system identification: a unified overview," Automatica, vol.31, issue 12, pp. 1691-1724, December 1995.
40. Specht, D.F., "A general regression neural network," Neural Networks, IEEE Transactions on, vol.2, no.6, pp.568-576, Nov 1991.
41. Strangas, E.G.; Aviyente, S.; Zaidi, S.S.H. , "Time-Frequency Analysis for Efficient Fault Diagnosis and Failure Prognosis for Interior Permanent-Magnet AC Motors," Industrial Electronics, IEEE Transactions on , vol.55, no.12, pp.4191-4199, Dec. 2008.
42. Sundararajan, D., A practical approach to Signals and Systems, Singapore: John Wiley & Sons, 2008, pp. 61-83.
43. Tobar, F.; Yacher, L.; Paredes, R.; Orchard, M., "Anomaly detection in power generation plants using similarity-based modeling and multivariate analysis," American Control Conference (ACC), 2011, pp.1940-1945, June 29 2011-July 1 2011.
44. Todorovic, V.M.; Tadic, P.R.; Djurovic, Z.M., "Expert system for fault detection and isolation of coal-shortage in thermal power plants," Control and Fault-Tolerant Systems (SysTol), 2010 Conference on, pp.666-671, 6-8 Oct. 2010.
45. Ye, H.; Wang, G.; Ding, S.X., "A new fault detection approach based on parity relation and stationary wavelet transform," American Control Conference, 2003. Proceedings of the 2003, vol.4, pp. 2991- 2996, 4-6 June 2003.