



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos



**Universidad de Valladolid**



**ESCUELA DE INGENIERÍAS  
INDUSTRIALES**

UNIVERSIDAD DE VALLADOLID

ESCUELA DE INGENIERIAS INDUSTRIALES

*Grado en Tecnologías Industriales*

**Control de calidad de un proceso mediante la detección y  
diagnóstico de anomalías usando técnicas de control  
estadístico de procesos**

Autor:

Pérez Franco, Iván

Tutor:

De La Fuente Aparicio,

María Jesús

Departamento de Ingeniería  
de Sistemas y Automática

Valladolid, junio 2020.



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos





## ÍNDICE DE CONTENIDOS

Resumen .....	7
Abstract.....	8
CAPITULO I: INTRODUCCIÓN .....	9
1.1. Motivación del proyecto .....	11
1.2. Objetivos .....	12
1.3. Organización de la memoria .....	14
CAPITULO II: MARCO TEÓRICO .....	17
2.1. Control de Calidad.....	19
2.1.1. Historia reciente.....	19
2.2. Control Estadístico de Procesos (SPC).....	21
2.2.1. Variabilidad en los procesos .....	22
2.2.2. Capacidad de un proceso .....	25
2.2.3. Gráficos de control.....	27
2.3. Análisis de Componentes Principales (PCA).....	30
2.3.1. Estadísticos empleados para monitorizar el proceso .....	33
2.3.2. Diagramas de contribución .....	35
2.4. Método t-SNE.....	36
2.4.1. Algoritmo t-SNE utilizado para la detección de fallos.....	40
2.5. Árboles de decisión.....	42
CAPITULO III: PLANTA TENNESSEE EASTMAN.....	47
3.1. Descripción del proceso .....	49
3.2. Datos del proceso .....	51
CAPITULO IV: TRABAJO REALIZADO .....	57
4.1. Metodología utilizada .....	59
4.2. Análisis de componentes principales (PCA) del proceso para detectar y diagnosticar los fallos producidos.....	60
4.2.1. Análisis del comportamiento normal del proceso .....	60
4.2.2. Detección de anomalías del proceso mediante PCA .....	63
4.2.3. Contribución de las variables al fallo .....	70
4.3. Uso del método t-SNE para la detección de fallos producidos en el sistema .....	75



4.3.1. Detección de anomalías del proceso usando la distancia euclídea .....	76
4.3.2. Detección de anomalías del proceso mediante reducción a 3 dimensiones con la técnica t-SNE.....	82
4.3.3. Detección de anomalías del proceso con la técnica t-SNE utilizando distintas distancias entre puntos.....	85
4.3.4. Mejora de la detección de anomalías del proceso utilizando la distancia Mahalanobis.....	91
4.4. Clasificación de las anomalías detectadas .....	95
4.4.1. Clasificación de fallos utilizando los datos de partida .....	96
4.4.2. Clasificación de fallos utilizando PCA.....	99
4.4.3. Clasificación de fallos utilizando t-SNE .....	101
CAPITULO V: CONCLUSIONES Y TRABAJO FUTURO .....	105
5.1. Comparación entre las técnicas de control estadístico de procesos utilizadas...	107
5.2. Trabajo que se podría realizar en un futuro .....	109
Bibliografía .....	111



## ÍNDICE DE FIGURAS

FIGURA 1: DISTINTAS DISTRIBUCIONES DE LA VARIACIÓN EN DOS PROCESOS [7] .....	23
FIGURA 2: ESQUEMA DE LAS 6 M'S PARA LA VARIABILIDAD DE UN PRODUCTO [8] .....	24
FIGURA 3: DISTINTOS PROCESOS EN FUNCIÓN DE SUS VALORES DE ESPECIFICACIÓN Y LA EXISTENCIA DE CAUSAS ESPECIALES [6] .....	25
FIGURA 4: RELACIÓN ENTRE LOS LÍMITES INFERIOR Y SUPERIOR CON LA DISTRIBUCIÓN DE UN PROCESO [7] .....	26
FIGURA 5: EJEMPLO DE UN GRÁFICO DE CONTROL PARA EL PESO DE UNA PIEZA [9] .....	28
FIGURA 6: VALORES DE LAS MUESTRAS TOMADAS Y REPRESENTACIÓN DE LOS SUBESPACIOS DE PCA POR LOS EJES DE UN ELIPSOIDE .....	31
FIGURA 7: ESPACIO X DE ALTA DIMENSIÓN EN 2D DONDE SE MIDEN SIMILITUDES POR PARES DE PUNTOS [13] .....	38
FIGURA 8: DISMINUCIÓN DE LA DIMENSIÓN DEL SISTEMA INICIAL X EN UNO DE MENOR DIMENSIÓN Y [14] .....	40
FIGURA 9: ESQUEMA DEL FUNCIONAMIENTO DE UN ÁRBOL DE DECISIÓN [16] .....	43
FIGURA 10: DIAGRAMA DE FLUJO DE TODO EL PROCESO DE LA PLANTA FICTICIA TENNESSEE EASTMAN [18] .....	50
FIGURA 11: COMPARACIÓN DEL ESTADÍSTICO $T^2$ PARA LOS DATOS DE COMPORTAMIENTO NORMAL CON EL UMBRAL $Ta2$ .....	62
FIGURA 12: COMPARACIÓN DEL ESTADÍSTICO Q (SPE) PARA LOS DATOS DE COMPORTAMIENTO NORMAL CON EL UMBRAL $Q_A$ .....	63
FIGURA 13: RESULTADOS DE LOS ESTADÍSTICOS $T^2$ (IZQUIERDA) Y Q (DERECHA) PARA EL FALLO 1 .....	64
FIGURA 14: RESULTADOS DE LOS ESTADÍSTICOS $T^2$ (IZQUIERDA) Y Q (DERECHA) PARA EL FALLO 8 .....	65
FIGURA 15: RESULTADOS DE LOS ESTADÍSTICOS $T^2$ (IZQUIERDA) Y Q (DERECHA) PARA EL FALLO 15 .....	66
FIGURA 16: CONTRIBUCIÓN DE CADA VARIABLE A QUE SE PRODUZCA EL FALLO 1 .....	71
FIGURA 17: CONTRIBUCIÓN DE CADA VARIABLE A QUE SE PRODUZCA EL FALLO 14 .....	72
FIGURA 18: COMPARACIÓN DE LOS VALORES CALCULADOS DEL ESTADÍSTICO $T^2$ CON EL UMBRAL $Ta2$ PARA UN COMPORTAMIENTO NORMAL DEL PROCESO .....	77
FIGURA 19: REPRESENTACIÓN DEL ESTADÍSTICO $T^2$ PARA LOS FALLOS 1 (IZQUIERDA) Y 8 (DERECHA) DEL SISTEMA MEDIANTE T-SNE .....	78
FIGURA 20: REPRESENTACIÓN DEL ESTADÍSTICO $T^2$ PARA EL FALLO 15 DEL SISTEMA MEDIANTE T-SNE .....	79
FIGURA 21: COMPARACIÓN DE LOS VALORES CALCULADOS DEL ESTADÍSTICO $T^2$ PARA UNA MATRIZ Y REDUCIDA A 3 DIMENSIONES CON EL UMBRAL $Ta2$ .....	83
FIGURA 22: GRÁFICO OBTENIDO MEDIANTE LA FUNCIÓN GSCATTER Y LA DISTANCIA EUCLÍDEA .....	86
FIGURA 23: GRÁFICO OBTENIDO MEDIANTE LA FUNCIÓN GSCATTER Y LA DISTANCIA MAHALANOBIS .....	87
FIGURA 24: GRÁFICO OBTENIDO MEDIANTE LA FUNCIÓN GSCATTER Y LA DISTANCIA COSENO .....	87
FIGURA 25: GRÁFICO OBTENIDO MEDIANTE LA FUNCIÓN GSCATTER Y LA DISTANCIA CHEBYCHEV .....	88



FIGURA 26: MATRIZ DE CONFUSIÓN UTILIZANDO LOS DATOS DE CREACIÓN DEL BOSQUE.....	97
FIGURA 27: MATRIZ DE CONFUSIÓN UTILIZANDO LOS DATOS DE TEST DEL BOSQUE .....	97
FIGURA 28: MATRIZ DE CONFUSIÓN MEDIANTE PCA UTILIZANDO LOS DATOS DE CREACIÓN DEL BOSQUE .....	100
FIGURA 29: MATRIZ DE CONFUSIÓN MEDIANTE PCA UTILIZANDO LOS DATOS DE TEST DEL BOSQUE .....	100
FIGURA 30: MATRIZ DE CONFUSIÓN MEDIANTE T-SNE UTILIZANDO LOS DATOS DE CREACIÓN DEL BOSQUE .....	102
FIGURA 31: MATRIZ DE CONFUSIÓN MEDIANTE T-SNE UTILIZANDO LOS DATOS DE TEST DEL BOSQUE.....	102



## Resumen

El término Industria 4.0 ha ganado mucha importancia a lo largo de estos últimos años debido a que cada vez son más las empresas que se nutren de la cantidad de información que aporta esta nueva revolución digital que ya muchos consideran como una candidata a ser la Cuarta Revolución Industrial. Uno de los pilares fundamentales que sustentan esta nueva tecnología es el concepto del *Big Data*, aportando una gran cantidad de datos que pueden ser tratados y estudiados de muchas maneras distintas.

En este trabajo gracias a esta tecnología se va a realizar un estudio para mejorar la calidad de un proceso industrial. Se utilizarán técnicas de control estadístico de procesos para la detección y diagnóstico de fallos o anomalías que puedan afectar negativamente al sistema y empeorar la calidad del producto final. Por un lado, se utiliza el Análisis de Componentes Principales (PCA), una técnica de reducción de dimensión lineal que nos permite trabajar con un menor número de variables a partir de las cuales se analiza el comportamiento de la planta. Por otra parte, se utiliza la técnica de incrustación de vecinos estocásticos distribuidos en  $t$  (t-SNE), que trabaja con una reducción de dimensionalidad no lineal. Estas técnicas se aplicarán a una planta química usada como benchmark en la literatura científica, la planta Tennessee Eastman (TEP), haciéndose un estudio comparativo entre ambos métodos.

**Palabras Clave:** Industria 4.0, Big Data, Planta Tennessee Eastman, Análisis de Componentes Principales (PCA), técnica de incrustación de vecinos estocásticos distribuidos en  $t$  (t-SNE), simulación, control estadístico de procesos, Matlab.



## Abstract

The term Industry 4.0 has become increasingly important over the last few years because numerous companies are relying on the amount of information provided by this new digital revolution that many consider it as a candidate to be the Fourth Industrial Revolution. One of the fundamental pillars that underpin this new technology is the concept of Big Data, providing a large amount of data that can be processed and studied in different ways.

In this project, thanks to this technology, a study will be carried out to improve the quality of an industrial process. Statistical process control techniques will be used to detect and diagnose failures or anomalies that can negatively affect the system and worsen the quality of the final product. On the one hand, Principal Component Analysis (PCA) is used, which is a linear dimension reduction technique that allows us to work with fewer variables that help to properly analyze plant simulation data. On the other hand, we use the T-distributed Stochastic Neighbor Embedding (t-SNE), later developed to PCA and working with a reduction of nonlinear dimensionality. These techniques will be applied to a chemical plant used as a benchmark in the scientific literature, the Tennessee Eastman plant (TEP), making a comparative study between both methods.

**Palabras Clave:** Industry 4.0, Big Data, Tennessee Eastman Plant, Principal Component Analysis (PCA), T-distributed Stochastic Neighbor Embedding (t-SNE), simulation, statistical process control. Matlab.





# CAPITULO I: INTRODUCCIÓN



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos





## 1.1. Motivación del proyecto

Todos los días se realizan procesos industriales que tienen un papel fundamental dentro del sector productivo debido al aumento en la demanda de producción originada por la calidad de vida a la que estamos actualmente acostumbrados gracias a los avances tecnológicos de los que continuamente se nutren estos procesos. No solo se ha conseguido reducir el tiempo de fabricación de la mayoría de los productos que consumimos, sino que también se proporciona una mayor calidad que logra satisfacer las actuales necesidades de los clientes. De esta manera, las empresas consiguen optimizar y mejorar el rendimiento de su proceso y los consumidores obtienen sus productos con mejores prestaciones y a un precio más económico.

En consecuencia, se ha aumentado de forma significativa la demanda por parte de la industria de sistemas de control que proporcionen una alta fiabilidad, que aseguren la calidad del producto final y la seguridad en los procesos de fabricación. Los efectos que generan los fallos en el proceso incrementan sus costes de operación y en ocasiones pueden llegar a resultar extremadamente peligrosos para la salud de las personas y generar un impacto medioambiental negativo en el entorno si no cumplen con todos los controles y medidas necesarias. Una advertencia anticipada puede ayudarnos a reaccionar antes de que se produzca una avería y evitar paradas innecesarias que supondrían un coste económico a la empresa e incluso evitar una catástrofe si el fallo fuera peligroso.

Para conseguir que esto funcione, se ha aumentado el nivel de automatización de los procesos industriales, creando sistemas de mayor complejidad a la hora de realizar las tareas de supervisión necesarias. Dichas tareas se realizan en gran medida por operadores humanos a través de los sistemas SCADA (Supervisory Control and Data Acquisition). Los operadores tienen que saber analizar la situación cuando alguna variable supera los



niveles de seguridad que se han estipulado, tomando las acciones correctivas para controlar el proceso. Si podemos ser capaces de implementar mecanismos de diagnóstico de fallos a estos sistemas de supervisión automática, conseguiremos que los operarios puedan controlar instalaciones muy complejas tomando las medidas correctivas necesarias en todo momento [1]. Los datos que se necesitan para poder controlar estos sistemas son más fáciles de obtener actualmente debido a la multitud de sensores que existen hoy en día y que nos permiten obtener la información de un gran número de variables de los procesos. Una vez obtenidos estos datos, la complejidad reside en relacionarlos entre sí de forma que nos aporten la información que necesitamos para diagnosticar posibles anomalías o fallos que ocurran en el proceso.

Las técnicas multivariantes surgen de esta necesidad, otorgando las herramientas necesarias al control estadístico de procesos y permitiendo la monitorización de estos procesos. El Análisis de Componentes Principales (PCA) es una de estas técnicas que consiste en reducir el número de variables originales a un número menor de variables, denominadas componentes principales, que serán combinación lineal de las variables iniciales y sintetizarán la mayor parte de información contenida en los datos originales. Otra de las técnicas basadas en la reducción de dimensiones utilizadas es la incrustación de vecinos estocástica distribuida (t-SNE), una técnica no lineal desarrollada en 2008 a diferencia de PCA que se desarrolló alrededor del año 1933.

## 1.2. Objetivos

En este trabajo se plantea mejorar la calidad de los procesos de producción, implementando métodos de detección y diagnóstico de fallos (FDD) o anomalías que pueden suceder en el sistema. Estos métodos que se van a desarrollar estarán basados en datos y serán suficientemente generales para



poder aplicarse a todo tipo de industrias, ya que en todas se pueden producir fallos que afectan a la fiabilidad del proceso, la seguridad de la planta, y a la calidad de los productos.

El primer paso de un método FDD es la detección del fallo o anomalía en la planta. Para resolver este problema se utilizarán técnicas de control estadístico multivariante de procesos, y en concreto se utilizarán dos técnicas, el análisis de componentes principales (PCA), acompañado de dos estadísticos, uno de ellos el estadístico Hotgelling's o  $T^2$  y el otro el estadístico Q o SPE. El segundo método que se va a utilizar es una técnica no lineal denominada t-SNE y trabajaremos también con su correspondiente estadístico  $T^2$ . Se hará una comparativa entre los resultados de detección de anomalías utilizando ambas herramientas.

Una vez detectado el fallo o la anomalía, es necesario diagnosticar que fallo ha ocurrido, o lo que es lo mismo, identificar el elemento o elementos de la planta que presentan problemas en su funcionamiento. En este caso se usará el diagrama de contribuciones a  $T^2$  y Q que nos dirán que variables son las que más afectadas están por los fallos, y por lo tanto que variables son las responsables del mismo. Esto indicará al operario donde centralizar su atención para que pueda reparar la anomalía, y que el sistema vuelva funcionar correctamente, y asegurar la calidad de los productos.

Otra segunda opción para resolver este problema de la detección y diagnóstico de fallos es plantear el problema como un problema de clasificación, en el que a partir de datos seamos capaces de clasificar si el sistema está en condiciones normales de funcionamiento o si hay un fallo, y en este caso nos indique que clase de fallo hay. Para resolver este problema usaremos una técnica muy actual de la inteligencia computacional llamada bosques aleatorios.



Todas estas técnicas, las aplicaremos sobre la planta Tennessee Eastman, un benchmark sacado de la literatura científica que sirve de referencia en el estudio del control y monitorización de procesos multivariantes.

### **1.3. Organización de la memoria**

Para organizar el presente trabajo, la memoria se ha estructurado en seis capítulos:

El primer capítulo corresponde con la introducción al tema del que trata este proyecto, poniendo en contexto y explicando los motivos por los cuales se ha propuesto realizar este trabajo. También se señalan los objetivos a alcanzar y la organización de la memoria.

En el segundo capítulo, se profundiza de forma teórica en el control de calidad, la detección de fallos y el análisis estadístico de procesos. También se explica cómo vamos a implementar las técnicas utilizadas para el análisis de los datos y su fundamento teórico.

En el tercer capítulo, se explica por qué se ha escogido la planta Tennessee Eastman como sujeto de pruebas y se realiza una explicación de su proceso industrial. Se identifican las variables y los fallos que se van a analizar.

En el cuarto capítulo, se documenta la parte práctica realizada, donde se explica detalladamente la metodología utilizada al aplicar cada una de las técnicas empleadas para la detección y análisis de los fallos producidos en la planta.

En el quinto capítulo, se presentan las conclusiones en base a los resultados obtenidos experimentalmente, analizando individual y globalmente las herramientas seleccionadas para este trabajo. Por otra parte, también se hacen unas sugerencias para futuros estudios relacionados con este tema.



Finalmente, el sexto capítulo de la memoria estará destinado a la recopilación de la bibliografía que se ha consultado.



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos







# CAPITULO II: MARCO TEÓRICO



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos





## 2.1. Control de Calidad

Actualmente, la mayoría de las empresas tienen contratados a jefes de departamento, directores de fábrica, encargados o supervisores que en conjunto se responsabilizan de la calidad de los productos o servicios que se ofrecen al consumidor, controlando que todo funcione correctamente en los respectivos departamentos, fábricas y equipos de empleados. A su vez, los técnicos especialistas e ingenieros se encargan de crear, implementar y mejorar las metodologías reflejadas en las normas que se utilizan en estos procesos sistemáticos, permitiendo a las empresas ofrecer la calidad de sus productos a un precio razonable.

El control de calidad presenta varios objetivos a parte de una mejora de las prestaciones del producto que quiere el cliente. El primero es consolidar la economía de un país ya que se logran exportar mayores dosis de productos de mayor calidad con un precio asequible, el segundo objetivo es asegurar unas bases económicas firmes que sustenten el futuro permitiendo establecer y exportar de forma continua la tecnología industrial. Finalmente, las empresas que mejoran y dan importancia a estos procesos de control, incrementan a su vez el nivel de vida de las personas, poniendo a disposición productos de buena calidad a un mayor número de clientes que se lo pueden permitir [2].

### 2.1.1. Historia reciente

Los sistemas de gestión y control de la calidad siempre han seguido un camino que prioriza la calidad y estandarización de los procesos industriales. Esto no es algo nuevo, sino que hace ya siglos los artesanos de la época se agrupaban en unas asociaciones económicas denominadas “gremios”, y con la llegada de la Primera Revolución Industrial comenzó el trabajo especializado y en serie dando lugar a la producción por procesos. Se



instauraron sistemas de gestión de calidad que implementaban ciertos estándares para el control de productos y resultados.

Tiempo después durante la Segunda Guerra Mundial creció la necesidad de poder simplificar procesos industriales sin renunciar a la seguridad, y por tanto, el estudio y mejora del control de calidad se convirtió en uno de los pilares más importantes sobre todo en la industria bélica dada la situación de la época. Con la ayuda de herramientas estadísticas como el muestreo y los gráficos de control, se crearon nuevas técnicas de control da calidad y muestreo de inspección y se anunciaron estándares y pautas de entrenamiento. Se creó el concepto de Calidad Total, definida por la ley ISO 9001:2015 como *“Estrategia de gestión de la organización que tiene como objetivo satisfacer de una manera equilibrada las necesidades y expectativas de todos sus grupos de interés, normalmente empleados, accionistas y la sociedad en general”*. [3]

Sin embargo, incluso finalizada esta guerra la importancia de la calidad siguió en aumento debido a la gran competitividad entre Japón y Estados Unidos. A finales del siglo XX nacieron los sistemas de gestión de calidad y posteriormente en el siglo XXI han surgido nuevas ideas como la responsabilidad social y la sostenibilidad que mejoran estos sistemas. [4].

Finalmente, llegamos al estado actual donde la gestión de calidad permite a las empresas alcanzar las especificaciones de sus clientes, incrementando la confianza del consumidor con la marca para posteriores pedidos. Además, si el consumidor se encuentra satisfecho recomendará la empresa a otras personas aumentando la red de clientes.

Por otra parte, también la gestión de calidad permite alcanzar los estándares internacionales reglamentarios en productos y servicios, lo que se transfiere en un uso eficiente de los recursos disponibles, aumentando las ganancias de la empresa y la confianza en los productos ofrecidos.



## 2.2. Control Estadístico de Procesos (SPC)

El Control Estadístico de procesos, generalmente conocido por SPC debido a sus siglas en inglés (*Statistical Process Control*), está formado por un conjunto de diversas técnicas y herramientas que tienen como finalidad examinar la calidad de un proceso productivo. Para ello, se controla a través de unos gráficos de control que identifican fácilmente las posibles anomalías que puedan perjudicar el proceso determinando si los resultados concuerdan con el diseño estipulado y si nos encontramos dentro de un rango de tolerancia previamente conocido. Si hay algún problema se alerta e informa a las personas encargadas en ese momento, que dispondrán de los datos necesarios para tomar las medidas correctivas correspondientes que permitan hacer funcionar al proceso de forma correcta. Si aplicamos con éxito el SPC se obtendrán grandes ventajas como por ejemplo la reducción de costes que provocan las averías de los productos, la disminución de piezas fuera de tolerancia o la reducción del tiempo de fabricación que permite cumplir con las fechas de entrega acordadas con los clientes.

La primera persona que desarrolló estos gráficos fue Walter Shewhart alrededor de 1920. Se encontraba en la empresa Western Electric dedicada a la fabricación de teléfonos. Dicha empresa estaba en auge en esa época debido al gran aumento de la demanda de sus productos, por lo tanto, tuvo que ampliar su plantilla de trabajadores e incrementó la cantidad de líneas de producción en gran medida lo que provocó que surgieran problemas de controlar el proceso correctamente. Posteriormente, los gráficos utilizados por Shewart fueron perfeccionados por Joseph Juran y Edwards Deming y junto a Walter Shewhart se les considera que son los precursores del control SPC que existe actualmente [5].



### 2.2.1. Variabilidad en los procesos

En todos los sistemas de producción, personas o máquinas trabajan con los materiales necesarios para la obtención de un producto con la calidad deseada, sin embargo, existen diversas fuentes de variación provocadas por la variabilidad inherente o natural cuyas causas no se pueden controlar, sin embargo, existe otra causa de diversidad debido a la variabilidad no natural la cual se puede controlar y en la que deberemos centrarnos para poder hacer una correcta gestión de la calidad de nuestro producto.

Por tanto, las causas de la variabilidad se pueden clasificar como:

- **Causas comunes:** Son causas aleatorias inherentes al proceso de fabricación con carácter permanente, son predecibles estadísticamente pero no se pueden evitar. No son determinantes para la pérdida de control del proceso, si solo existieran estas causas comunes el proceso estaría bajo control estadístico.
- **Causas especiales:** Son variaciones no inherentes al proceso con carácter esporádico o puntual. Su existencia en el proceso implica que este abandone su estado de control por lo que es necesario que los responsables del control de calidad apliquen las medidas necesarias para identificarlas y si es posible eliminarlas [6].

Walter Shewhart como ya mencionamos anteriormente fue uno de los pioneros del control de calidad y fue la primera persona que se percató que los datos que provienen de las mediciones de las variaciones de un proceso productivo la mayoría de las veces no producen un gráfico de distribución normal sino que tienden a organizarse en torno a una distribución distinta de la campana de Gauss. Como se puede observar en la figura 1, la variación de un proceso puede tener distintas distribuciones, siendo primordial reducir esa

variabilidad para no salirnos de los límites de calidad establecidos y no generar productos defectuosos en cuanto a calidad se refiere.

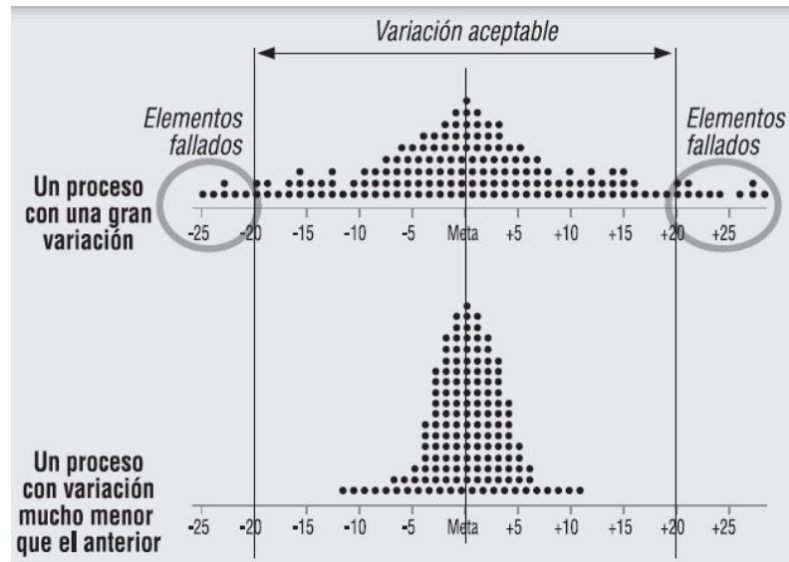


Figura 1: Distintas distribuciones de la variación en dos procesos [7]

Para poder determinar estadísticamente esta variabilidad se realiza un experimento aleatorio, donde se elige una muestra al azar y se mide la magnitud de la variable deseada que intervenga en las especificaciones de calidad de nuestro producto. Los factores que no podemos controlar implicarán que obtengamos resultados diferentes por lo que la magnitud que estamos controlando se puede interpretar como una variable aleatoria que no solo estará afectada por la variabilidad de la muestra sino que afectará también la variabilidad a la que está sometido el instrumento de medida utilizado en el muestreo. Por tanto, la variabilidad detectada en el experimento será la suma de la variabilidad real del producto debido al material y todo lo relacionado con el proceso de fabricación utilizado, más la variabilidad del sistema de medición [8]. Un esquema que puede resultar de utilidad para recordar los distintos factores que afectan a la variabilidad de un



proceso es el esquema de las 6 M's (Figura 2).

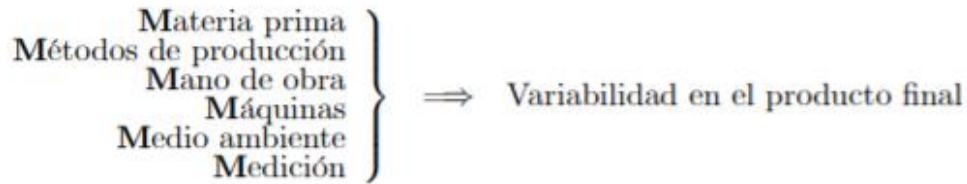


Figura 2: Esquema de las 6 M's para la variabilidad de un producto [8]

Además, en el estudio de la variabilidad de un proceso deberemos tener en cuenta las especificaciones de calidad de nuestro producto para saber si estamos en un rango de variabilidad adecuado. Para ello, se definen los límites superior e inferior dentro de los cuales se cumplen las especificaciones. Hay que discernir entre los límites de las especificaciones del cliente, denominados LSL el límite inferior y USL el límite superior, y los límites naturales que siguen una distribución normal del propio proceso denominados LPL el límite inferior y UPL el límite superior. Si ambos límites de las especificaciones en valor absoluto son mayores que los límites naturales del proceso, es posible afirmar que el proceso es estable y capaz por lo que el producto obtenido tendrá una alta probabilidad de estar en el rango de las características deseadas logrando una alta fiabilidad del proceso de producción. Sin embargo, el sistema puede ser estable pero no capaz al mismo tiempo debido a que uno de los límites de las especificaciones sea menor en valor absoluto que el correspondiente límite natural. En consecuencia, no se podrá asegurar que el proceso productivo nos otorgue un producto con la calidad deseada. Finalmente, si ocurren las causas especiales mencionadas anteriormente se producirá un fallo en el proceso que implicará una inestabilidad que hará que el sistema se encuentre fuera de control estadístico y que hay que evitar en todo proceso. En la Figura 3 se pueden observar los distintos casos mencionados.



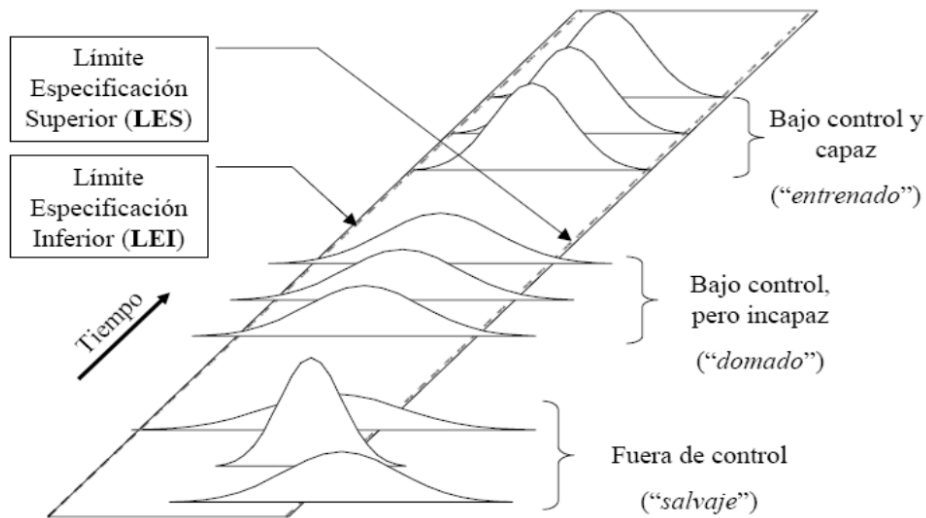


Figura 3: Distintos procesos en función de sus valores de especificación y la existencia de causas especiales [6]

### 2.2.2. Capacidad de un proceso

El control estadístico de procesos ayuda a realizar y mantener la media y varianza del proceso constantes y así lograr mantener la distribución necesaria para cumplir con la calidad final del producto. Sin embargo, que tengamos un buen control estadístico del proceso no significa que este genere los productos adecuados a las especificaciones de diseño ya que los límites establecidos para este control se basan en la media y variabilidad de la distribución de muestreo y no en las especificaciones técnicas del diseño.

Debido a esto, se busca que el proceso tenga la capacidad para cumplir de forma adecuada estas especificaciones. En la figura 4(a), podemos observar que el proceso es capaz ya que los valores extremos de la distribución se encuentran dentro de los límites superior e inferior de las especificaciones. Por otra parte, en la figura 4(b), el proceso no tiene la capacidad de producir la mayoría de los productos con la calidad requerida, ya que habrá un porcentaje de piezas que queden por debajo del límite inferior deseado.

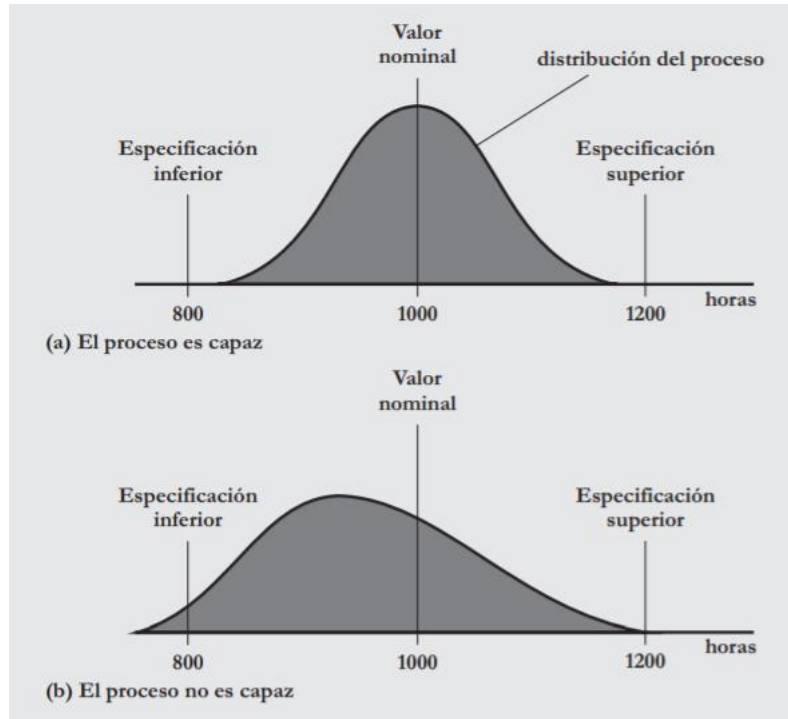


Figura 4: Relación entre los límites inferior y superior con la distribución de un proceso [7]

Otro aspecto importante destacable es la razón de capacidad de proceso  $C_{pk}$ , definida por la siguiente fórmula:

$$C_{pk} = \min \left[ \frac{\bar{x} - \text{Especificación inferior}}{3\sigma} ; \frac{\text{Especificación superior} - \bar{x}}{3\sigma} \right] \quad (1)$$

Se ha establecido que la mayoría de los valores provenientes de la distribución de un proceso se encuentran en un rango de  $\pm 3\sigma$  siendo  $\sigma$  la desviación estándar de la media por lo que el rango total es de seis desviaciones estándar para que un valor medido tenga la calidad que buscamos. Si el proceso que estamos controlando queremos que sea capaz deberemos tener que la amplitud de tolerancia (diferencia entre las especificaciones superior e inferior) sea mayor que  $6\sigma$ . Si la razón de capacidad  $C_{pk}$  es mayor que 1.0, como por ejemplo 1.54, el proceso será



capaz. Sin embargo, si el  $C_{pk}$  es menor que 1.0 se estarán generando bastantes productos fuera del intervalo de tolerancia deseado.

Por otro lado, también se define el índice de capacidad de proceso  $C_p$  con el objetivo de comparar la capacidad del proceso y la amplitud del intervalo de tolerancia establecido.

$$C_p = \frac{\text{Especificación superior} - \text{Especificación inferior}}{6\sigma} \quad (2)$$

Si deseamos que nuestro sistema de producción nos otorgue productos que estén dentro de este rango de tolerancia, el valor de  $C_p$  tendrá que ser mayor que 1 [8].

### 2.2.3. Gráficos de control

Los gráficos de control son una herramienta ampliamente utilizada en el control estadístico de procesos que permite observar de forma gráfica el sistema que estamos analizando y así poder detectar de forma sencilla las posibles anomalías. Como se puede observar en la figura 5, este gráfico se representa mediante puntos que representan el valor de la magnitud medida a lo largo del tiempo y se unen mediante líneas que permiten analizarlo rápidamente con vistazo. También se grafican los límites de control inferior y superior y la línea central que representa los valores que deberían tomar las muestras si no hubiera variabilidad en el proceso. Esta es la característica más importante de estos gráficos ya que los hace muy intuitivos y fácilmente interpretables.

Otra característica de estos gráficos es la estimación de parámetros como la media o desviación típica del sistema, que nos permite saber si el proceso tiene la capacidad necesaria para cumplir con las especificaciones. Además, proporcionan los datos necesarios para conocer las variaciones que se



producen en un momento determinado, evitando que se ajusten los procesos de forma innecesaria si dichas variaciones son simplemente ruido o en caso contrario, que las causas de estas variaciones sean por una causa asignable. El hecho de poder conocer las causas de las variaciones otorga una mejora de la productividad y eficiencia del análisis de estos datos.



Figura 5: Ejemplo de un gráfico de control para el peso de una pieza [9]

Los primeros gráficos de control utilizados fueron los del ingeniero Shewhart, quien realizó a grandes rasgos los primeros controles que se basaban en el estudio de la media del proceso y del rango o desviación dependiendo la situación. Se realizaban solo estudios univariantes, donde solo se medía y controlaba una variable del sistema. Los gráficos que representan la evolución de la media de las muestras de tamaño  $n$  se denominan gráficos para  $\bar{X}$ . Por otro lado, están los gráficos del rango que tienen en cuenta la diferencia  $R$  entre el mayor valor y el menor que puede tomar una muestra de un conjunto  $n$ .

Los límites aplicados a este tipo de gráficos son los siguientes:

$$LCS = \mu + 3 \frac{\sigma}{\sqrt{n}} \quad (3)$$

$$LC = \mu \quad (4)$$



$$LCI = \mu - 3 \frac{\sigma}{\sqrt{n}} \quad (5)$$

El valor tres sigma se suele tomar para establecer los límites de control ya que corresponde a un valor de 0.0027 del parámetro estadístico  $\alpha$  para distribuciones normales. A este parámetro  $\alpha$  se le denomina nivel de significación, es un concepto estadístico que otorga el valor de la probabilidad de obtener falsas alarmas. Dependiendo lo estricto que sea el control que estamos realizando se puede variar el valor de este parámetro, teniendo en cuenta que cuanto mayor sea tendremos más falsas alarmas y el control no será muy preciso. Además, los datos de la media  $\mu$  y de la desviación típica  $\sigma$  no suelen ser conocidos y es necesario estimarlos estadísticamente.

Por otra parte, una alternativa a estos gráficos propuestos por Shewhart son los gráficos CUSUM, los cuales tienen en cuenta la información de muestras anteriores permitiendo observar rápidamente si existen pequeñas variaciones en las muestras e identificar visualmente los cambios en la media del proceso. En estos gráficos se representan la suma acumulada de los datos de las desviaciones típicas de las muestras respecto a su media estadística. Funciona de forma más eficiente que los gráficos de Shewhart si la desviación respecto a la media es menor que dos veces  $\sigma$ . Sin embargo, su inconveniente es que son lentos a la hora de detectar grandes cambios en el sistema debido a que los datos representados tienen ese carácter acumulativo y al existir una correlación entre ellos se hace difícil poder interpretar el patrón que siguen [10]



## 2.3. Análisis de Componentes Principales (PCA)

En el análisis de datos multivariantes uno de los mayores problemas es la cantidad de variables con las que hay que trabajar, todas las variables medibles del proceso, por lo que es necesario implementar un método de reducción de la dimensionalidad; y de esta manera, reduciremos los datos que serán analizados a cambio de una pequeña pérdida de información.

Si tenemos  $n$  muestras de  $m$  variables del sistema, el objetivo del Análisis de Componentes Principales es analizar si podemos representar el proceso con un menor número de variables construidas a raíz de una combinación lineal de las variables originales. De esta manera se consigue representar el sistema en un espacio de dimensión reducida con pequeñas observaciones del espacio general que involucra a todas las variables y tiene dimensión  $m$ .

Esta herramienta se basa en una técnica de proyección que transforma ortogonalmente las variables del sistema inicial, las cuales generalmente están correlacionadas entre sí, en conjunto de variables linealmente no correlacionadas y de menor dimensión que las anteriores. Estas nuevas variables son las que se denominan componentes principales.

Como podemos observar en la figura 6, el conjunto de datos de  $m$  dimensiones se ajusta a un elipsoide que tiene dimensión  $p$  de tal manera que la dimensión  $p$  sea menor que  $m$ . PCA nos permite dividir el espacio en dos subespacios diferentes, cada eje del elipsoide representa uno de esos subespacios y aporta cierta información del sistema. El conjunto de ejes menores representado por el eje verde de la figura 6 representa una varianza baja a lo largo del mismo por lo que si este eje es sustituido por su correspondiente componente principal, se pierde simplemente una pequeña cantidad de información. Por otro lado, el eje rojo captura la tendencia del proceso. Se busca de esta manera el menor número de componentes principales que pueda capturar la máxima variabilidad de los datos originales.

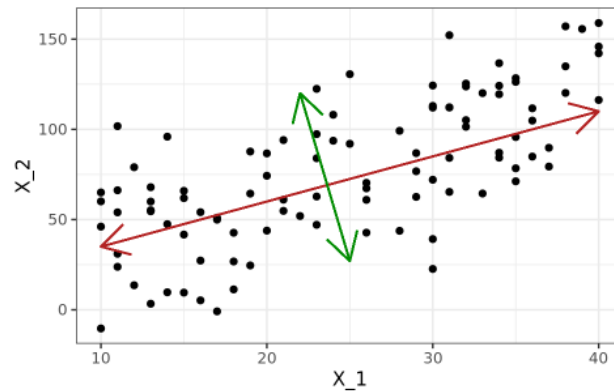


Figura 6: Valores de las muestras tomadas y representación de los subespacios de PCA por los ejes de un elipsoide

Para aplicar el método PCA es necesario realizar un pre-tratamiento de los datos del proceso, que consiste en eliminar las variables inadecuadas como pueden ser las que tengan errores muy grandes debido a los instrumentos de medida y que por tanto no aportan unos datos reales del proceso. También tenemos que realizar un escalado para asegurar que todas las variables tienen la misma influencia en el proceso de monitorización, para ello, se intenta capturar la variación de la media restando a cada variable su valor medio y se busca por otra parte tener varianza unitaria dividiendo cada variable por su desviación estándar.

Una vez hemos realizado el pre-tratamiento de los datos, comenzamos por aplicar el método tomando como ejemplo una matriz  $X$  de datos de dimensiones  $n \times m$ , tal que  $X \in \mathbb{R}^{n \times m}$ , siendo  $n$  el número de observaciones realizadas y  $m$  las variables del sistema.

$$X = \begin{pmatrix} X_{11} & \cdots & X_{1m} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{nm} \end{pmatrix} \quad (6)$$



El objetivo es determinar un conjunto de vectores propios ortogonales que se ordenen según la cantidad de variabilidad que aportan al sistema. Para ello, utilizamos la matriz X (6) para hallar la matriz de correlación R (7):

$$R = \frac{1}{n-1} X^T X \quad (7)$$

A continuación, se halla la descomposición en valores singulares de R de esta manera:

$$R = V \cdot \Lambda \cdot V^T \quad (8)$$

Siendo  $\Lambda$  una matriz que dispone de los valores propios positivos de R y ordenados de mayor a menor de tal manera que  $\lambda_1 \geq \lambda_2 \dots \geq \lambda_m \geq 0$ . Cabe destacar que cada uno de estos valores propios es igual a su varianza:  $\lambda_i = \sigma_i^2$ . Por otro lado, la matriz V es una matriz ortogonal cuyas columnas son los vectores propios (o vectores de carga) de R, conocidos normalmente en PCA como *scores*.

Para poder representar la máxima variabilidad de los datos sin que afecte el ruido, se escoge un número  $a$  que corresponde con los vectores de carga de mayor magnitud. Se crea una matriz P que contiene estos vectores de dimensión  $m \times a$ , tal que  $P \in \mathfrak{R}^{m \times a}$ , con esta matriz P y la matriz de datos X (6) podemos proyectar los datos iniciales en un espacio reducido creando la matriz de transformación T, tal que  $T \in \mathfrak{R}^{n \times a}$ . Se denomina *scores* a cada valor que compone esta nueva matriz.

$$T = X \cdot P \quad (9)$$

Si se definen cada columna  $i$ -ésima de T como  $t_i$  obtenemos las siguientes características:

- $var(t_1) \geq var(t_2) \geq \dots \geq var(t_a)$  es decir la varianza se encuentra ordenada de mayor a menor.





- Se encuentra centrado en la media:  $\mu(t_i) = 0 \forall i$
- Presenta descomposición ortogonal:  $t_i^T \cdot t_j = 0; \forall i \neq j$
- No existe otra expansión ortogonal que contenga  $a$  componentes y tenga más información de la variación de los datos analizados.

Habiendo hecho la transformación y obtenido la matriz T se puede recalcularse la matriz original X de la siguiente manera:

$$\hat{X} = T \cdot P^T \quad (10)$$

Con esta matriz recalculada podemos obtener la matriz de residuos E a partir de la diferencia entre la matriz original X (6) y la matriz calculada posteriormente  $\hat{X}$  (10):

$$E = X - \hat{X} \quad (11)$$

De esta manera, con el espacio de residuos E obtenemos la variabilidad de los datos de observación a través de los vectores propios que están asociados a los valores singulares  $m-a$  de menor valor y nos permite representar al conjunto de datos de partida:

$$X = \hat{X} + E \quad (12)$$

### 2.3.1. Estadísticos empleados para monitorizar el proceso

Para la detección de fallos con PCA se van a utilizar dos estadísticas con el objetivo de determinar un umbral que permita detectar cualquier anomalía en las muestras obtenidas y así notificar si ha ocurrido un fallo en el proceso.

La primera que utilizaremos será la **Estadística de Hotelling ( $T^2$ )**, partiendo de la matriz X (6) se obtiene una observación  $x \in \mathfrak{R}^{m \times 1}$  y utilizando la matriz P de



componentes principales podemos calcular esta estadística de la siguiente manera:

$$T^2 = x^T \cdot P \cdot \Lambda_a^{-1} \cdot P^T \cdot x \quad (13)$$

Se utiliza la inversa de  $\Lambda_a$  que está formado por las primeras  $a$  filas y columnas, teniendo en cuenta que  $a$  es la cantidad de componentes principales que hemos escogido.  $T^2$  se puede interpretar como la distancia de la observación a la media del modelo si se calcula para una muestra  $x$  de  $m$  variables.

Una vez hemos calculado la matriz  $T^2$  deberemos comparar cada dato con un umbral  $T_a^2$ , de tal manera que si una muestra supera el umbral  $T^2 > T_a^2$  detectaremos una alarma en el proceso que nos avisará si ha ocurrido un fallo en el sistema y el proceso está fuera de control estadístico. El umbral se calculará a través de la siguiente fórmula:

$$T_a^2 = \frac{(n^2 - 1)a}{n(n - a)} F_\alpha(a, n - a) \quad (14)$$

Para hallar el umbral se utiliza la distribución de Fisher-Snedecor  $F_\alpha(a, n - a)$ , donde  $\alpha$  es el valor del nivel de significancia que queremos que tenga nuestro sistema de detección, representando el grado de compromiso de las falsas alarmas que normalmente tiene un rango de valores entre 0.05 y 0.01.

Para monitorizar los restantes  $m-a$  variables vamos a utilizar la **Estadística Q** o también conocida como **SPE**. Esta estadística se puede calcular a través del vector de residuos  $r$ :

$$r = (I - P \cdot P^T) \cdot x \quad (15)$$

Con ayuda de este vector, se define la estadística Q de la siguiente manera:



$$Q = r^T \cdot r \quad (16)$$

A continuación, necesitaremos calcular el umbral para poder evaluar los datos obtenidos como hemos hecho con la estadística  $T^2$  anteriormente. El umbral  $Q_a$  se obtiene mediante la siguiente fórmula:

$$Q_a = \theta_1 \left[ \frac{h_0 c_{\alpha} \sqrt{2\theta_2}}{\theta_1} + 1 + \frac{\theta_2 h_0 (h_0 - 1)}{\theta_1^2} \right]^{\frac{1}{h_0}} \quad (17)$$

Siendo  $h_0 = 1 - \frac{2\theta_1\theta_3}{3\theta_2^2}$  y  $\theta_i = \sum_{j=a+1}^m \lambda_j^i$

### 2.3.2. Diagramas de contribución

Podemos estudiar qué variables originales del proceso son las responsables de un fallo detectado utilizando los diagramas de contribución para ambas estadísticas  $Q$  y  $T^2$ .

En primer lugar, para la estadística  $T^2$  hay que determinar los  $r$  scores  $t_i$  responsables de que el sistema se encuentre fuera de control, tal que  $r < a$  siendo  $a$  el número de componentes principales determinado tras aplicar el método de PCA. Para ello, calculamos la matriz  $T$  como se hizo en la ecuación 10, pero solo en el instante que ha ocurrido el fallo:

$$T = X(t_{fallo}, :) \cdot P \quad (18)$$

Elevamos esta matriz al cuadrado para obtener todos los datos positivos y dibujamos un diagrama de barras en el que observamos cuales son las componentes principales que se encuentran fuera de control.

A continuación, se calcula la contribución de cada variable del sistema original a que la aparición de anomalías en las componentes principales:



$$cont_{ij} = \frac{t_i}{\lambda_i} p_{ij} x_j \quad (19)$$

Finalmente, se haya la suma de todas las contribuciones de cada variable y se representa un diagrama de barras para observar fácilmente cual tiene mayor contribución y así detectar las causas que han producido el fallo del sistema.

$$CONT_j = \sum_{i=1}^r cont_{ij} \quad (20)$$

Por otro lado, el cálculo de la contribución mediante el estadístico Q resulta mucho más sencillo. Simplemente, se trata de calcular el vector de residuos  $r$  (Ecuación 15) y lo elevamos al cuadrado para no tener valores negativos como hemos hecho en la ecuación 18 para el estadístico  $T^2$ . Finalmente, solo tenemos que realizar un diagrama de barras de  $r^2$  y obtendremos la contribución de todas las variables a que se haya producido el fallo en el sistema.

## 2.4. Método t-SNE

El siguiente método de reducción dimensional que se va a utilizar es la técnica de incrustación de vecinos estocásticos distribuidos en  $t$ , comúnmente denominado t-SNE (*t-distributed stochastic neighbor embedding*), este método fue creado por Laurens van der Maaten y Geoffrey Hinton y realiza una reducción no lineal de los datos de partida disminuyéndolos hasta un espacio de 2 o 3 dimensiones [11]

En primer lugar, este método se basa en la técnica *Stochastic Neighbor Embedding* (SNE) desarrollada por Hinton y Roweis, la cual transforma un



espacio euclídeo donde se encuentran las distancias que existen entre un gran número de datos  $X$  de alta dimensión, a unas probabilidades condicionadas que reflejan las similitudes que hay entre ellos. En el espacio  $X$ , la similitud entre  $x_i$  y  $x_j$  se representa por la probabilidad  $p_{ij}$  condicionada a que  $x_i$  fuera elegida por  $x_j$  como vecino. El espacio al que queremos reducir este conjunto de datos será  $Y$ , donde existen otras probabilidades condicionadas  $q_{ij}$  para los puntos  $y_i$  e  $y_j$  que se obtienen a partir de los respectivos puntos  $x_i$  y  $x_j$  del conjunto de datos principal. Si el método se realiza correctamente, se podrá comprobar que las probabilidades condicionales de los puntos  $x_i$  y  $x_j$  será la misma que la de los puntos  $y_i$  e  $y_j$  de la dimensión  $Y$  reducida. De esta manera, SNE busca minimizar la diferencia entre  $p_{ij}$  y  $q_{ij}$  a través de una representación de baja dimensión.

Sin embargo, si se utiliza el método t-SNE se puede llegar a conseguir una mejor visualización y optimización debido a que disminuye las posibilidades de que las observaciones se encuentren todas juntas en el medio del mapa de visualización.

De la misma manera que ocurre con el método SNE, el primer paso de t-SNE es también medir las similitudes que hay entre los diversos pares de puntos de un espacio  $X$  de alta dimensión como podemos observar en la figura 7.

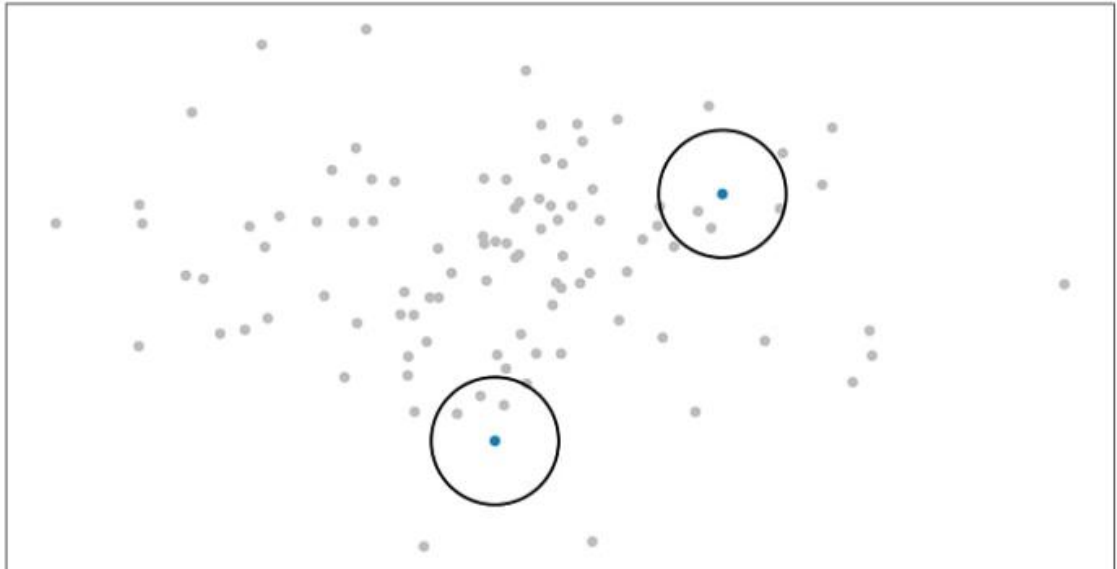


Figura 7: Espacio  $X$  de alta dimensión en 2D donde se miden similitudes por pares de puntos [13]

En el espacio  $X$ , se toma un punto  $x_i$  y la distribución gaussiana centrada en ese punto, reflejada por los círculos que aparecen en la figura 7. Se mide la densidad de todos los puntos que hay dentro del área que encierra la distribución gaussiana. Se normalizan los demás puntos y se calcula la probabilidad que refleja de manera proporcional la similitud entre  $x_i$  y  $x_j$ , obteniendo un conjunto de probabilidades  $P_{ij}$  para todos los puntos a través de la ecuación 21.

$$p_{ij} = \frac{\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)}{\sum_k \sum_{l \neq k} \exp\left(-\frac{\|x_k - x_l\|^2}{2\sigma^2}\right)} \quad (21)$$



Para los puntos que se encuentran a una distancia pequeña, la probabilidad de elegir un par de puntos que presenten similitudes es muy grande, por otro lado, si los puntos se encuentran a grandes distancias la probabilidad será muy pequeña.

El siguiente paso se centra en el espacio reducido  $Y$ , de dos o tres dimensiones, al que se desea llegar. Se realiza un procedimiento similar al anterior, buscando un segundo conjunto de probabilidades  $Q_{ij}$  para el espacio  $Y$ . Se representa cada punto  $x_i \in X$  en este nuevo espacio tal que  $y_i \in Y$ . A continuación, en lugar de utilizar la distribución gaussiana para definir un área como hicimos en el paso anterior, se utiliza una distribución en  $t$  de Student con un solo grado de libertad (también denominada distribución de Cauchy). De esta manera, se consigue mejorar la visualización ya que conseguimos que dos puntos que se encuentre a gran distancia en el espacio inicial  $X$ , se encuentren todavía más lejanos en el espacio de baja dimensión  $Y$ . Por último, se normaliza dividiendo entre todos los pares de puntos para medir la densidad entre  $y_i$  e  $y_j$ , consiguiendo la similitud  $q_{ij}$  entre pares de puntos en  $Y$  a través de la ecuación 22:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|y_k - y_l\|^2)^{-1}} \quad (22)$$

Se busca que la diferencia entre  $q_{ij}$  y  $p_{ij}$  sea la menor posible para que la estructura que aparece en el mapa de visualización de los datos sea muy parecida y así poder trabajar con el espacio  $Y$  de dimensión más baja de manera similar a si estuviéramos utilizando el espacio  $X$  de dimensión mucho mayor. Finalmente se mide la diferencia entre las distribuciones de probabilidad de los espacios bidimensionales  $q_{ij}$  y  $p_{ij}$  utilizando un método estándar denominado divergencia de Kullback-Leibler que se usa para medir distancias entre este tipo de distribuciones de probabilidad. De esta forma, la

técnica t-SNE ordena los puntos minimizando  $KL(P || Q)$  para conseguir que los valores de  $p_{ij}$  sean parecidos a  $q_{ij}$ . [12]

$$KL(P||Q) = \sum_i \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (23)$$

### 2.4.1. Algoritmo t-SNE utilizado para la detección de fallos

Como observamos en la figura 8, los datos del espacio Y se calculan mediante un método de manifold learning (en este caso el método t-SNE) a partir del espacio de alta dimensión de datos X. Pero, también puede considerarse que existe una matriz de proyección lineal A, desde el espacio de alta dimensión al de baja dimensión) que permite transformar los datos iniciales en un espacio de dos o tres dimensiones.

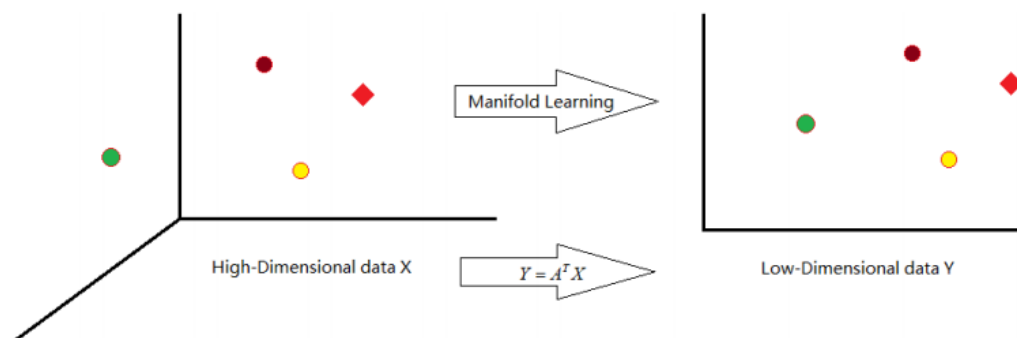


Figura 8: Disminución de la dimensión del sistema inicial X en uno de menor dimensión Y [14]

De esta manera, sea el espacio original de entrenamiento  $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{n \times m}$  con n observaciones de un proceso industrial, de tal manera que  $x_i \in \mathbb{R}^m$  ( $i = 1, 2, \dots, n$ ) donde cada observación está formada a su vez por m variables distintas medidas a través de sensores en un instante i de tiempo. Por lo tanto, obtenemos un vector de dimensión m para cada muestra  $x_i$ . Por





otro lado tenemos el espacio reducido  $Y = [y_1, y_1, \dots, y_n] \mathbb{R}^{l \times n}$  donde  $y_i \in \mathbb{R}^l$  ( $i = 1, 2, \dots, n$ ) son los datos proyectados de  $x_i$  en un espacio dimensional  $L$  siendo  $L$  la dimensión reducida que normalmente tomará el valor de 2 o 3 calculados con el método t-SNE

$$x_i \rightarrow y_i = A^T x_i \quad (24)$$

Basándose en la hipótesis de que el mapeo t-SNE se realice localmente, se establece una proyección lineal para aproximar el mapeo que relaciona el espacio de alta dimensionalidad y el espacio de baja dimensionalidad embebido; teniendo en cuenta que  $A$  es una matriz  $m \times l$ , tal que  $A \in \mathbb{R}^{m \times l}$  y puede ser obtenida a través de la siguiente regresión lineal de mínimos cuadrados de los datos de entrenamiento en los que no existen fallos en el sistema.

$$A = (XX^T)^{-1}XY^T \quad (25)$$

El siguiente paso es calcular la estadística  $T^2$ . La estadística  $T^2$  calculada para PCA se basa en la distancia de Mahalanobis, aquí ahora usaremos la técnica t-SNE que mapea los datos del proceso desde la dimensión alta a la baja dimensión preservando la estructura local de los mismos [14]. Se podría calcular la matriz de dimensión  $Y$  reducida a través de la matriz  $A$  calculada con los datos de entrenamiento de funcionamiento normal de la planta.

$$Y = A^T X \quad (26)$$

Finalmente calculamos el estadístico  $T^2$  de los datos de comportamiento normal, siendo  $y_i$  la correspondiente proyección de cada dato  $x_i$  del proceso a través de la siguiente ecuación:

$$T_i^2 = y_i^T \left( \frac{YY^T}{n-1} \right)^{-1} y_i \quad (27)$$



A continuación se elige un umbral  $T_a^2$  para el control estadístico del sistema, de forma que si el valor de  $T^2$  es superior al valor del umbral  $T_a^2$  significaría que el sistema está detectando una anomalía en el sistema.

Para comprobar si el método funciona, se calcula  $T^2$  para unos nuevos datos que al contrario que los anteriores, estarán caracterizados por un fallo del sistema y el sistema deberá detectar anomalías en el proceso. Se calculan los nuevos datos  $y_{nuevo}$  de la dimensión reducida  $Y$  a través de los nuevos datos  $x_{nuevo}$  y la matriz  $A$  calculada anteriormente con los datos de comportamiento normal.

$$y_{nuevo} = A^T x_{nuevo} \quad (28)$$

Una vez calculamos todos los datos con la ecuación 28, se procede a calcular la estadística  $T_{nuevo}^2$  para cada nuevo dato a través de la siguiente fórmula:

$$T_{nuevo}^2 = y_{nuevo}^T \left( \frac{Y Y^T}{n-1} \right)^{-1} y_{nuevo} \quad (29)$$

Finalmente, solo queda comparar los datos  $T_{nuevo}^2$  calculados con el umbral anterior  $T_a^2$  de tal manera que si  $T_{nuevo}^2 < T_a^2$  significa que el dato representa un comportamiento normal del proceso en ese instante, en caso contrario habrá una anomalía correspondiente a un fallo en el proceso [14].

## 2.5. Árboles de decisión

En el sector del aprendizaje automático una de las herramientas más utilizadas son los árboles de decisión. Un árbol de decisión es una estructura ramificada que muestra las consecuencias de diferentes opciones y nos permite clasificar los distintos casos que tenemos. Existen puntos donde hay

que tomar decisiones denominados nodos de decisión, estos están unidos por las ramas partiendo de un nodo principal para finalmente llegar a las hojas que representan las últimas decisiones. De esta manera, si queremos clasificar un nuevo caso, se compararán los atributos con las decisiones que se toman en los nodos y se abrirá camino por la rama que coincida con los nodos que posean dichos valores. Por último, se llega a la hoja que predice la clase que corresponde con el caso tratado.

El árbol se basa en las diferentes características de los datos de entrenamiento para aprender los factores necesarios para inferir las etiquetas de clase de los ejemplos. Para desglosar el conjunto de datos que se analiza, los árboles de decisión utilizan un algoritmo que formula preguntas hasta recoger la suficiente información que les permita hacer una predicción. A partir del nodo raíz donde comienza el algoritmo, se dividen los datos en la característica que aporte la mayor información posible hasta que se llega a las hojas, todo ello de forma iterativa [15].

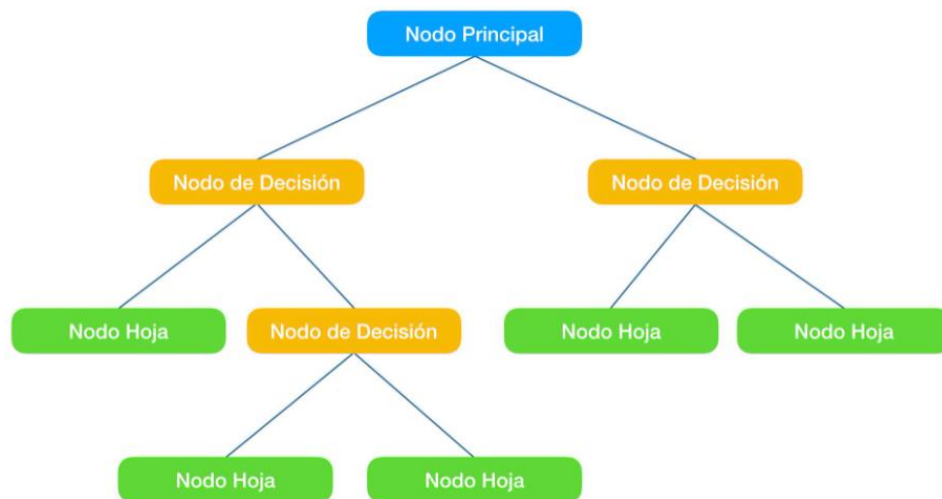


Figura 9: Esquema del funcionamiento de un árbol de decisión [16]



Existen diversos parámetros que caracterizan estos modelos de clasificación. El primer de ellos es la máxima profundidad, que se atribuye a la mayor longitud que hay desde el nudo inicial hasta las hojas. A mayor profundidad puede generarse un sobreajuste que se puede evitar estableciendo un parámetro con la máxima longitud. Por el contrario, a poca profundidad se producirá un sub-ajuste.

El siguiente parámetro a tener en cuenta es el máximo número de muestras que influenciará en el total de nodos que tendrá el árbol. Si descartamos uno de los nodos puede ocurrir que la mayoría de las muestras se encuentre en la parte eliminada lo que se traduce en un uso ineficiente de los recursos. Sin embargo, esto se puede evitar estableciendo en cada hoja un número máximo de muestras que puedan albergar. El número máximo es uno de los parámetros más importantes ya que en la mayoría de las ocasiones tendremos muchos datos para construir el árbol y cada vez que cortemos un nudo tendremos que revisar para cada característica todo el conjunto de datos, lo que provoca que se hagan demasiadas operaciones. Esto se puede evitar limitando el número de características que buscamos al realizar cada corte, teniendo en cuenta que si este número es lo bastante grande, será muy probable que encontremos la característica deseada entre aquellas que buscamos. No obstante, si este número no es tan alto como el total de características causará un aumento considerable en la velocidad de los cálculos.

En este trabajo vamos a utilizar un tipo de árboles de decisión denominados *Random Forest* (Bosques aleatorios) ya que es una técnica muy común cuando vamos a analizar datos con muchas características como sucede en nuestro proceso en el que intervienen bastantes variables. Random forest es una técnica basada en la combinación de distintos árboles predictivos denominados clasificadores débiles. Trabaja con una colección de árboles no correlacionados y los promedia, cada árbol depende de un vector aleatorio de



la muestra de forma independiente y que presenta la misma distribución de todos los árboles del bosque. De esta manera, se crea un algoritmo que primero proyecte una muestra de inicio de tamaño  $n$ . Después, a partir de la muestra principal se diseña un árbol de decisión donde cada nodo tendrá características aleatoriamente seleccionadas sin reemplazamiento.

Este proceso se repetirá un número  $x$  de veces y se añadirá la predicción seleccionada para cada árbol, asignado por mayoría la etiqueta del tipo de clase.



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos





# **CAPITULO III: PLANTA TENNESSEE EASTMAN**



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos







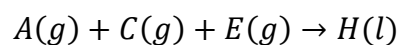
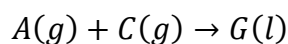
### 3.1. Descripción del proceso

En este trabajo se va a utilizar los datos extraídos del proceso Tennessee Eastman debido a que ha sido muy utilizado por muchos científicos e ingenieros para experimentar en el ámbito del control estadístico de procesos.

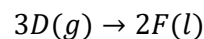
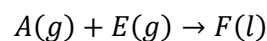
Todo comenzó a finales del siglo XX cuando la compañía Eastman Chemical sacó adelante en 1993 un proyecto junto a la universidad de Tennessee publicado por J. J. Downs y E. F. Vogel. Consistía en crear un modelo experimental basado en una planta real de la propia compañía química. Se trata de un proceso caracterizado por su alta no linealidad y complejidad, que está formado por un gran número de variables distintas que permiten crear multitud de escenarios distintos, permitiendo poder estudiar y analizar distintos métodos de detección y diagnóstico de fallos de la ingeniería de control de procesos [18].

Como ya se ha comentado, el proceso trata de una planta química en la que se dan lugar cuatro reacciones con un total de ocho componentes los cuales son: un inerte (B), cuatro reactivos (A, C, D, E), dos productos (G, H) y un subproducto (F). Las reacciones químicas que tienen lugar son las siguientes:

- Dos reacciones para la formación de los productos G y H:



- Dos reacciones para la formación del subproducto F:



Todas las reacciones están caracterizadas por ser reversibles, exotérmicas, el equilibrio es función de la temperatura y sus velocidades se pueden asemejar a una cinética de reacción de orden uno [19].

En la figura 10 se observa el diagrama de flujo de todo el proceso que se lleva a cabo en la planta. En primer lugar, se observa que existen cinco operaciones unitarias: una de reacción de los componentes en el reactor, una de condensación, una de compresión, una de separación vapor-líquido y finalmente una operación de desorción o también llamada stripping. Por otra parte, también hay otros elementos secundarios como válvulas, indicadores, cromatógrafos de gases y demás elementos de control.

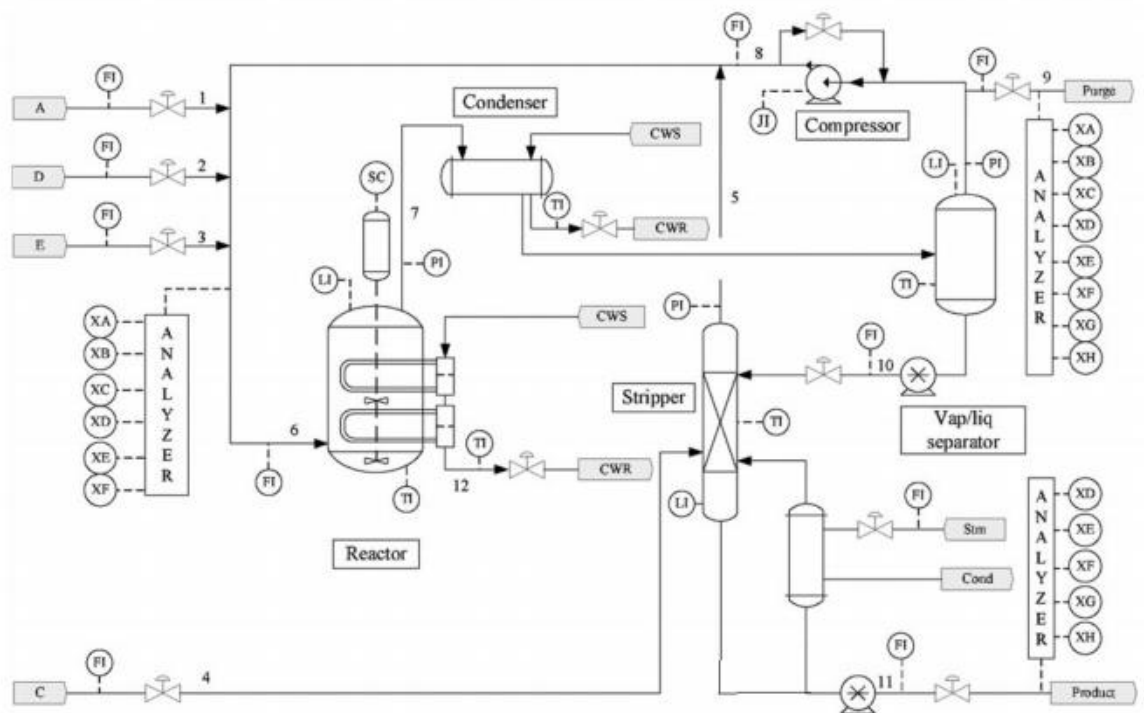


Figura 10: Diagrama de flujo de todo el proceso de la planta ficticia Tennessee Eastman [18]

Al comienzo del proceso se introducen los reactivos gaseosos A, D, C y E al reactor, teniendo en cuenta que el reactivo C pasa previamente por una operación de desorción antes de llegar a juntarse con los demás productos en la corriente de entrada al reactor. Además, se añade un disolvente catalizador en fase líquida que acelera el proceso de reacción y no abandona el reactor al ser no volátil. También se extrae el calor producido por las reacciones



exotérmicas a través de un intercambiador de calor situado en el reactor que permite disminuir su temperatura.

Las salidas del reactor son los productos en estado vapor aunque también se encuentran reactivos que no han reaccionado y que tendrán que recircularse para aumentar la eficacia del proceso. Todo ello se hace pasar por un condensador para enfriarlos y hacer que pasen a fase líquida, seguidamente, se introducen a un separador vapor-líquido donde los componentes que después de esta operación se encuentran en fase gaseosa por un lado se purgan y por otro se hacen pasar por un compresor que los recircula a la entrada del reactor. Por otra parte, los componentes condensados se redirigen a una columna de destilación alimentada por el reactivo C y que permite eliminar por desorción los reactivos condensados en la anterior operación. Finalmente, el producto de colas de la operación de stripping está formado por los componentes inerte B y el subproducto F que son purgados del proceso al hacerlo pasar por un separador vapor-líquido y también se encuentran los productos G y H deseados que son refinados en una sección no incluida al no tener importancia para el estudio que se va a realizar.

### 3.2. Datos del proceso

A continuación se va a presentar todos los datos con los que vamos a trabajar para poder analizar los distintos métodos de detección y clasificación de fallos de los que hemos hablado previamente. El modelo de planta Tennessee Eastman cuenta con un total de 12 variables manipuladas (XMV) y 41 variables medidas (XMEAS) presentadas en las tablas 1 y 2 respectivamente. Cabe destacar que de las 41 variables medidas, las 22 primeras se miden de forma continua mientras que las 19 restantes son obtenidas de los distintos analizadores que se encuentran en la planta que dan valores cada un intervalo discreto de tiempo. Además, las variables manipuladas (XMV) permiten al diseñador tener hasta un total de 12 grados de libertad que



residen en manipular 9 válvulas de flujo, 2 válvulas de control de temperatura y una velocidad de agitación. Para poder clasificar las anomalías que ocurren en el proceso también se tienen los datos de los tipos de fallos que ocurren a lo largo de toda la planta química y se pueden encontrar en la tabla 3 mostrada a continuación con las anteriores tablas mencionadas [19]:

Número de variable	Descripción de variable	Unidades
XMV (1)	Flujo de alimentación D	$\text{Kgh}^{-1}$
XMV (2)	Flujo de alimentación E	$\text{Kgh}^{-1}$
XMV (3)	Flujo de alimentación A	$\text{kscmh}$
XMV (4)	Flujo de alimentación A y C	$\text{kscmh}$
XMV (5)	Válvula de recirculación del compresor	%
XMV (6)	Válvula de purga	%
XMV (7)	Flujo de líquido del separador LV	$\text{m}^3\text{h}^{-1}$
XMV (8)	Flujo de líquido de la columna de stripping	$\text{m}^3\text{h}^{-1}$
XMV (9)	Válvula de vapor de la columna de stripping	%
XMV (10)	Flujo de agua de refrigeración del reactor	$\text{m}^3\text{h}^{-1}$
XMV (11)	Flujo de agua en el condensador	$\text{m}^3\text{h}^{-1}$
XMV (12)	Velocidad del agitador del reactor	$\text{rpm}$

Tabla 1: Variables manipuladas de la planta química [18]



Número de variable	Descripción de la variable	Unidades
XMEAS (1)	Flujo de alimentación de A	kscmh
XMEAS (2)	Flujo de alimentación de A	kscmh
XMEAS (3)	Flujo de alimentación de A	kscmh
XMEAS (4)	Flujo de alimentación de A y C	kscmh
XMEAS (5)	Flujo de recirculación	kscmh
XMEAS (6)	Flujo de alimentación al reactor	kscmh
XMEAS (7)	Presión del reactor	kPa
XMEAS (8)	Nivel del reactor	%
XMEAS (9)	Temperatura del reactor	°C
XMEAS (10)	Flujo de purga	Kscmh
XMEAS (11)	Temperatura del separador	°C
XMEAS (12)	Nivel del separador	%
XMEAS (13)	Presión del separador	kPa
XMEAS (14)	Corriente del separador	m <sup>3</sup> h <sup>-1</sup>
XMEAS (15)	Nivel del destilador (stripper)	%
XMEAS (16)	Presión del destilador (stripper)	kPa
XMEAS (17)	Corriente del destilador (stripper)	m <sup>3</sup> h <sup>-1</sup>
XMEAS (18)	Temperatura del destilador (stripper)	°C
XMEAS (19)	Flujo de vapor del destilador (stripper)	Kgh <sup>-1</sup>
XMEAS (20)	Potencia de compresor	kW
XMEAS (21)	Temperatura de la salida de agua de refrigeración del reactor	°C
XMEAS (22)	Temperatura de la salida de agua de refrigeración del separador	°C
XMEAS (23-28)	Concentración de la alimentación del reactor (A-F)	% mol
XMEAS (29-36)	Concentración de la purga (A-H)	% mol
XMEAS (37-41)	Concentración aguas abajo del destilador (A-H)	% mol

Tabla 2: Variables medidas de la planta química [18]



Tipo de fallo	Descripción del fallo	Unidades
IDV (1)	Relación de flujo de alimentaciones A/C, composición de B constante	Escalón
IDV (2)	Composición de B con relación A/C constante	Escalón
IDV (3)	Temperatura de alimentación D	Escalón
IDV (4)	Temperatura de entrada del agua al refrigerante del reactor	Escalón
IDV (5)	Temperatura de entrada del agua al refrigerante del condensador	Escalón
IDV (6)	Pérdida de alimentación de A	Escalón
IDV (7)	Pérdida de presión en la corriente C	Escalón
IDV (8)	Composición de las alimentaciones A, B y C	Variación aleatoria
IDV (9)	Temperatura de alimentación D	Variación aleatoria
IDV (10)	Temperatura de alimentación C	Variación aleatoria
IDV (11)	Temperatura de entrada del agua refrigerante al reactor	Variación aleatoria
IDV (12)	Temperatura de entrada del agua al refrigerante del condensador	Variación aleatoria
IDV (13)	Cinética de reacciones	Variación lenta
IDV (14)	Válvula del agua de refrigerante del reactor	Bloqueo
IDV (15)	Válvula del agua de refrigerante del condensador	Bloqueo
IDV (16)	Desconocido	Desconocido
IDV (17)	Desconocido	Desconocido
IDV (18)	Desconocido	Desconocido
IDV (19)	Desconocido	Desconocido
IDV (20)	Desconocido	Desconocido
IDV (21)	Desconocido	Constante

Tabla 3: Clasificación de los tipos de fallo de la planta química [18]



Los datos del modelo teórico de la planta Tennessee Eastman son públicos y perfectamente visibles para todo el mundo. Para este trabajo, se han utilizado diferentes conjuntos de datos. El primero consta de una simulación de la planta ante un comportamiento normal sin ningún tipo de fallo. Después, se añaden otros 21 ensayos realizados cada uno de ellos para un tipo de fallo diferente. Cabe destacar que los fallos 3,9 y 15 se les denomina fallos incipientes debido a que son muy difíciles de detectar por los métodos

Estos datos pueden ser obtenidos fácilmente a través del siguiente enlace (<http://web.mit.edu/braatzgroup/links.html>).



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos







# **CAPITULO IV: TRABAJO REALIZADO**



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos





## 4.1. Metodología utilizada

Como ya hemos comentado anteriormente, en este trabajo se va a tratar de analizar los datos de simulación del proceso Tennessee Eastman de manera que se puedan detectar y diagnosticar las anomalías durante su funcionamiento, a través de las técnicas de análisis multivariante ya explicadas en el capítulo II con el objetivo de experimentar y observar cómo funcionan estos métodos para poder finalmente obtener un buen control de calidad del proceso.

Se han utilizado dos bloques de ficheros que contienen distintos datos de simulación del proceso. Cada bloque está constituido por 22 simulaciones, la primera simulación contiene los datos de funcionamiento normal de la planta sin que ocurra ningún tipo de fallo. Las 21 simulaciones restantes corresponden respectivamente a cada uno de los 21 fallos que se han catalogado y que podemos encontrar en la tabla 3 del capítulo anterior.

El primer bloque está constituido por 21 simulaciones de 480 observaciones, es decir 480 instantes de tiempo en los que se miden todas las variables del proceso. Además, la simulación de datos de comportamiento normal del proceso consta de 500 observaciones en vez de 480. Por otra parte, el segundo bloque de datos que se va a utilizar está compuesto por simulaciones que cuentan cada una de ellas con un total de 960 observaciones. Sin embargo, a diferencia del anterior bloque, en este las primeras 160 observaciones siempre registrarán datos de comportamiento normal previo al fallo, a excepción de la simulación de funcionamiento normal que constará de 960 observaciones sin fallo. Por tanto, para las 21 simulaciones con fallo habrá que distinguir dos tramos: El primero desde el instante de observación 1 hasta el 160, y el segundo desde el instante de medición 161 hasta el 960 donde ya está presente el fallo correspondiente a la simulación con la que se esté trabajando. De esta manera, podemos



comprobar si la técnica que estamos utilizando para la detección de anomalías funciona correctamente en distintas ocasiones.

En relación con los datos simulados, cabe mencionar a los tres siguientes fallos: IDV (3) y IDV (9) relacionados con la temperatura de alimentación de D y IDV (15) relacionado con la válvula del agua de refrigerante del condensador. Son denominados fallos incipientes debido a su gran dificultad para ser detectados por lo que no deberemos preocuparnos si los métodos que estamos utilizando no logran detectarlos de forma correcta.

Por último hay que destacar que en este trabajo se va a utilizar el software MATLAB™ en su versión R2021a para analizar los datos de simulación y utilizar las técnicas de control estadístico de procesos.

## **4.2. Análisis de componentes principales (PCA) del proceso para detectar y diagnosticar los fallos producidos**

La primera técnica de análisis multivariante que vamos a utilizar para intentar detectar los fallos que se producen en el sistema es el análisis de componentes principales. Primero analizaremos los datos de comportamiento normal reduciendo el número de variables a las componentes principales que reduzcan el espacio de trabajo. También calcularemos los estadísticos ya comentados y sus correspondientes umbrales para después poder trabajar con los datos de simulación con fallo en la planta.

### **4.2.1. Análisis del comportamiento normal del proceso**

En primer lugar, se importan los datos de comportamiento normal a través de una matriz  $X$  de dimensión  $500 \times 52$  por lo que nuestro espacio muestral principal de datos es  $\mathfrak{R}^{500 \times 52}$  el cual vamos a tratar de reducir a un espacio de  $\mathfrak{R}^{500 \times a}$  siendo  $a$  el número de componentes principales que vamos a obtener al



aplicar este método. Esta matriz  $X$  de datos la tenemos que normalizar a media cero y varianza uno para que todos los datos tengan la misma relevancia y funcione de forma correcta el método de componentes principales. Se obtiene la matriz de correlación  $R$  a partir de la matriz  $X$  (Ecuación 7) y extraemos de la matriz de correlación sus valores singulares de manera decreciente, es decir, de mayor a menor valor numérico. Estos valores singulares se guardan en la matriz  $\Lambda$  y también se crea la matriz  $V$  compuesta por los vectores propios (vectores de carga) de  $R$  (Ecuación 8). El siguiente paso es elegir con qué variabilidad vamos a trabajar, en este caso se ha elegido una variabilidad del 90%. Una vez elegido este parámetro, se utiliza un test de porcentaje de varianza para calcular los componentes principales, de manera que se van sumando en orden decreciente los valores singulares de la matriz  $\Lambda$  hasta que se llega al porcentaje de varianza límite que en nuestro caso es 90%.

Después de realizar el test de varianza, se obtiene un total de 31 componentes principales y se define la matriz  $P \in \mathfrak{R}^{52 \times 31}$  que contiene los vectores de carga correspondiente a los valores singulares de las 31 variables con los valores más altos que acumulan el nivel de variabilidad aceptado. Con la ayuda de esta nueva matriz  $P$  podemos proyectar el espacio inicial  $\mathfrak{R}^{500 \times 52}$  a una proyección reducida, que denominaremos  $T$  (Ecuación 9) y que tiene una dimensión  $\mathfrak{R}^{500 \times 31}$ .

Para finalizar con los datos de comportamiento normal, se calculan los estadísticos  $T^2$  y  $Q$  que serán la forma que tendremos de visualizar si hay alguna anomalía en el proceso. Primero se calcula el estadístico  $T^2$  a través de la ecuación 13 para lo que es necesario elegir un nivel de significancia  $\alpha$  que en nuestro caso hemos elegido 0.02. Obtendremos una matriz de  $1 \times 500$ , es decir, un valor del estadístico  $T^2$  para cada observación. Para comprobar si este valor obtenido corresponde a un comportamiento normal del proceso, se calcula el umbral  $T_\alpha^2$  (Ecuación 14) que se obtendrá gracias a los datos con



los que estamos trabajando de comportamiento normal de la planta y que será utilizado en las siguientes simulaciones con fallo de planta para tratar de detectar anomalías. Aplicando estos pasos el valor del umbral  $T_a^2$  que hemos obtenido es 53,6117.

En la figura 11 se observan en azul los datos correspondientes al estadístico  $T^2$  que acabamos de calcular para cada una de las 500 observaciones. Se comprueba que no se supera el umbral  $T_a^2$  calculado menos en un instante, sin embargo, este valor es despreciable ya que podemos observar que todos los demás se encuentran por debajo del umbral y el proceso está funcionando correctamente.

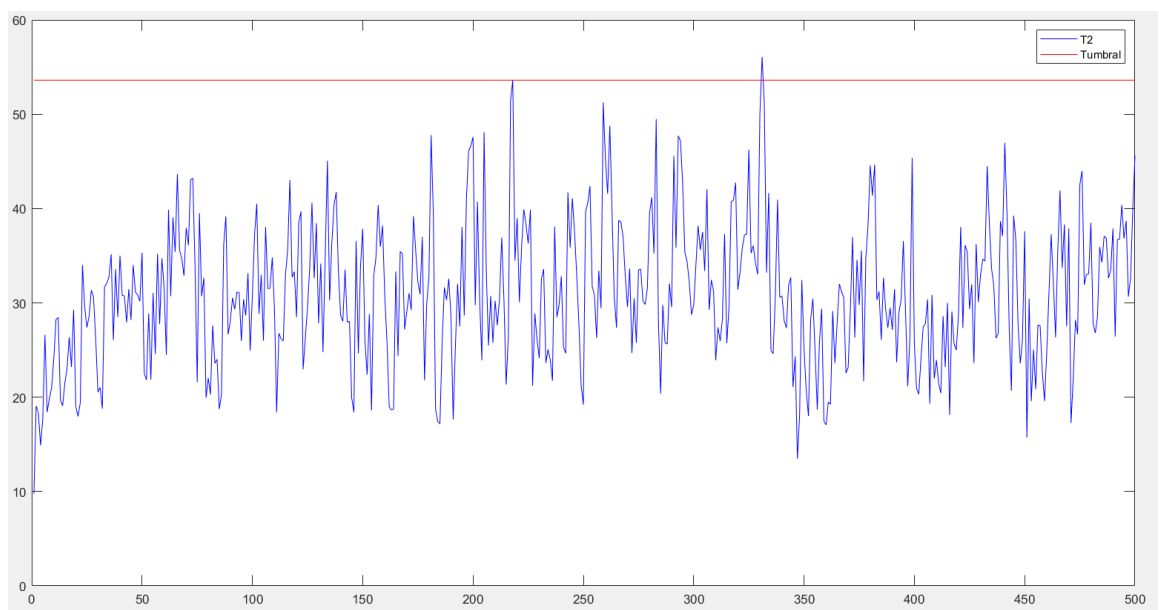


Figura 11: Comparación del estadístico  $T^2$  para los datos de comportamiento normal con el umbral  $T_a^2$

Por otro lado, se calcula también el estadístico Q (Ecuación 16) a través del vector  $r$  de residuos calculado con la ecuación 15. Siguiendo la misma metodología, se calcula el umbral  $Q_a$  (Ecuación 17) para el cual obtenemos un valor de 10,5719 que nos servirá para las demás simulaciones de la planta y que nos permitirá averiguar si ha habido algún fallo en el sistema. Al igual que

en la figura anterior correspondiente a  $T^2$ , en la figura 12 se observa cómo solo hay 4 instantes aislados donde el valor de  $Q$  supera el umbral de  $Q_a$

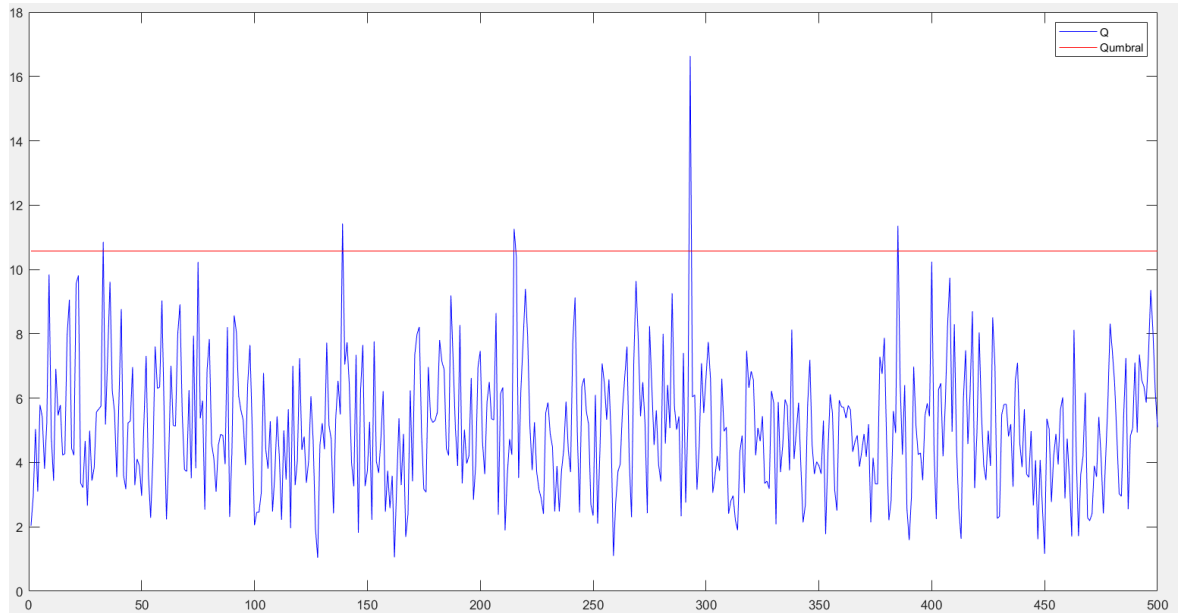


Figura 12: Comparación del estadístico  $Q$  (SPE) para los datos de comportamiento normal con el umbral  $Q_a$

#### 4.2.2. Detección de anomalías del proceso mediante PCA

El siguiente paso después de haber calculado el modelo PCA para el funcionamiento normal de la planta es utilizar este modelo para detectar los posibles fallos que pueden ocurrir en el sistema. Para ello, se utilizan las 21 simulaciones restantes caracterizadas cada una de ellas por un fallo de la planta. Al igual que en el apartado anterior, normalizamos la matriz de datos inicial  $X_i$  a una matriz  $Y_i$  de media cero y varianza y dimensión  $480 \times 52$  correspondiente a las 480 instantes de observación y las 52 variables del sistema.

Ahora simplemente tendremos que calcular los estadísticos  $T^2$  y  $Q$  y compararlos con sus correspondientes umbrales  $T_a^2$  y  $Q_a$ . Para evitar falsas alarmas, el método que se va a utilizar para averiguar si realmente ha



ocurrido un fallo en el sistema es que será necesario que los valores de  $T^2$  pasen el umbral  $T_a^2$  10 veces de forma consecutiva.

En primer lugar, calculamos el estadístico  $T^2$  (Ecuación 9) utilizando la matriz  $P$  y los valores singulares de mayor valor correspondientes a las componentes principales calculada anteriormente para los datos de comportamiento normal. Después se calcula el estadístico  $Q$  (Ecuación 16) a través de la matriz de residuos  $r$  calculada con la matriz  $P$  del apartado anterior y la nueva matriz  $Y_i$  de datos actuales.

A continuación, se han comparado todos los estadísticos de las 21 simulaciones correspondientes a ambos bloques de datos, los que tienen 480 observaciones y los que tienen 960. Sin embargo, como son muchos datos se van a utilizar figuras de solo unos pocos fallos que muestren el resultado del uso del método PCA. Para ello, se ha seleccionado el fallo 1, el fallo 8 y el fallo 15. Este último fallo comprobaremos que es muy difícil de detectar como ya comentamos anteriormente.

Finalmente, se mostrará una tabla que recoja todos los datos necesarios de todas las simulaciones de la planta que tenemos y así poder analizar en conjunto el resultado de aplicar esta técnica para detectar anomalías en el proceso.

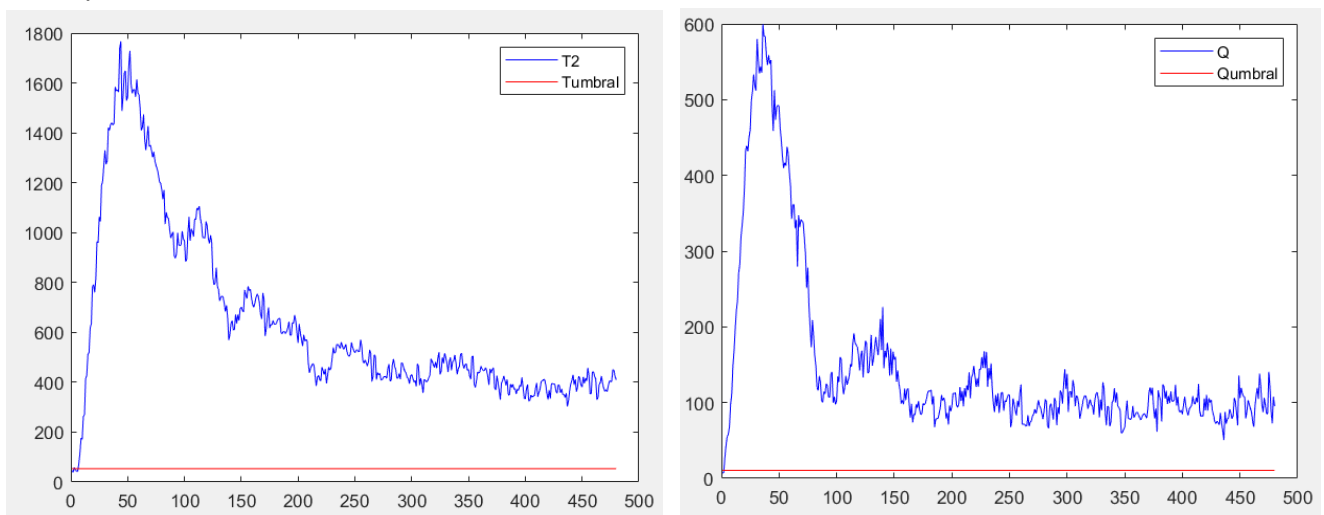


Figura 13: Resultados de los estadísticos  $T^2$  (izquierda) y  $Q$  (derecha) para el fallo 1



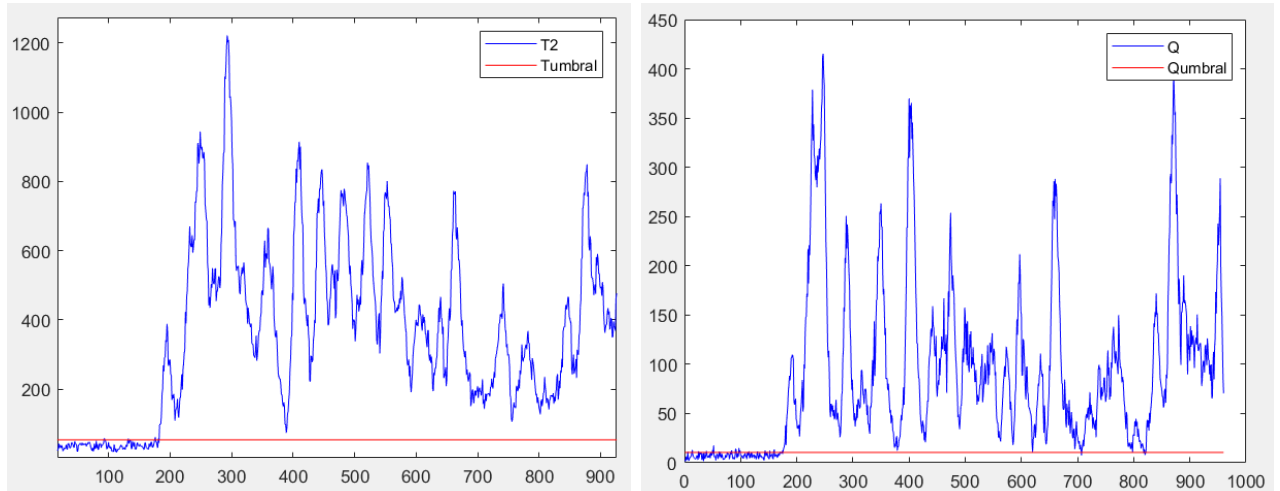


Figura 14: Resultados de los estadísticos T2 (izquierda) y Q (derecha) para el fallo 8

En la figura 13, se observan los resultados de aplicar el método PCA para detectar el Fallo 1 “Relación de flujo de alimentaciones A/C, composición de B constante “. Se ha utilizado el bloque de datos con 480 observaciones donde existe fallo en el sistema desde el primer instante. La línea azul corresponde con el estadístico calculado para cada instante y la línea roja el umbral que si es superado 10 veces consecutivas se puede afirmar que el método detecta que ha habido fallo. Comprobamos que para ambos estadísticos  $T^2$  y Q, se supera el umbral correspondiente desde prácticamente los primeros instantes por lo que el método detecta rápidamente el fallo que se ha producido en el proceso. También se observa que ambos estadísticos siguen una evolución muy similar a lo largo de las observaciones debido a que en el momento que ocurre el fallo el sistema se desestabiliza en gran medida pero después consigue regularse en torno a un valor menor que sigue estando fuera del umbral establecido.

Por otra parte, en la figura 14 se han utilizado el bloque de datos que contienen 960 observaciones para el fallo 8 “Composición de las alimentaciones A, B y C”. Podemos comprobar que el método detecta perfectamente que las primeras 160 muestras corresponden a datos de

comportamiento ya que ambos estadísticos se encuentran por debajo del umbral. A partir de la muestra 161 comienza a ocurrir el fallo en el sistema y observamos como rápidamente los valores de los estadísticos suben por encima del umbral. Al contrario que en el caso anterior, el valor de los estadísticos no se estabiliza en torno a un valor sino que observamos como la evolución de los datos tanto de  $T^2$  como de  $Q$  varían cada cierto intervalo de tiempo.

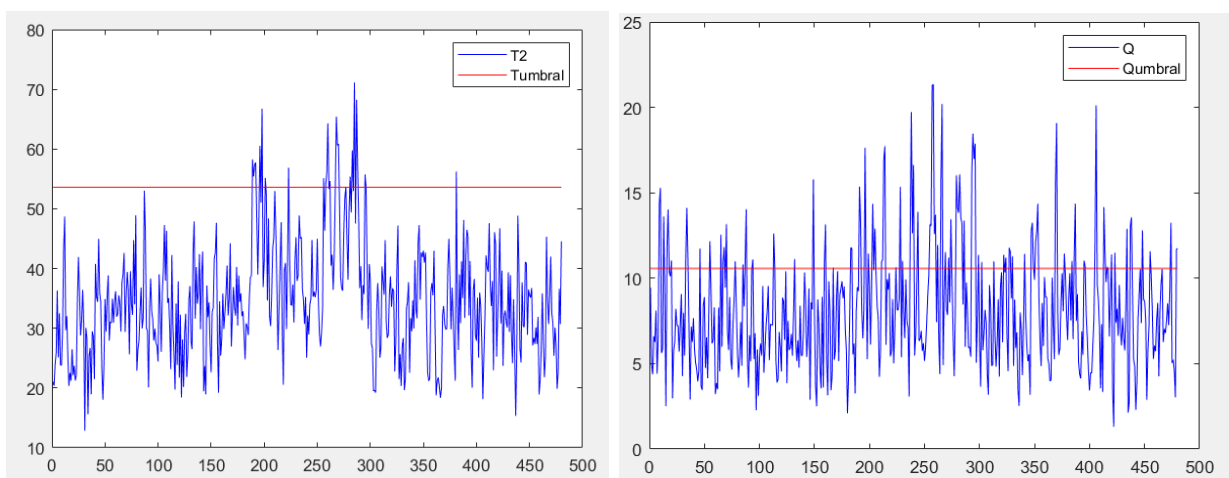


Figura 15: Resultados de los estadísticos  $T^2$  (izquierda) y  $Q$  (derecha) para el fallo 15

En último lugar, en la figura 15 se muestran los datos en relación con el fallo 15 “Válvula del agua de refrigerante del condensador”. Este fallo sabemos que es uno de los más difíciles de detectar, esto deriva en que pocos valores de los estadísticos  $T^2$  y  $Q$  calculados superan el umbral necesario para que el método PCA detecte que hay una anomalía en el sistema. En ninguno de los dos estadísticos se supera el umbral 10 veces consecutivas por lo que este método no logra detectar el fallo 15 del proceso.

Se ha recopilado todos los datos que se han obtenido al analizar todos las simulaciones del proceso ante todos los posibles fallos que pueden ocurrir en la planta Tennessee Eastman. Se ha calculado el tiempo que se tarda en



detectar el fallo y el número de alarmas detectadas, es decir, el número de veces que se supera el umbral establecido.

### Resultados obtenidos utilizando el estadístico T<sup>2</sup>:

Tipo de fallo	Instante en el que se detecta el fallo	Alarmas detectadas (%)
IDV (1)	7	98,96
IDV (2)	12	97,71
IDV (3)	No detecta	2,08
IDV (4)	No detecta	56,88
IDV (5)	1	42,50
IDV (7)	1	100
IDV (8)	17	96,67
IDV (9)	No detecta	2,50
IDV (10)	104	48,54
IDV (11)	68	56,88
IDV (12)	25	95,42
IDV (13)	21	95,83
IDV (14)	11	99,38
IDV (15)	No detecta	5,00
IDV (16)	No detecta	6,04
IDV (17)	52	78,75
IDV (18)	93	81,04
IDV (19)	No detecta	15,00
IDV (20)	39	50,42
IDV (21)	452	13,13
IDV_te (1)	7	99,38
IDV_te (2)	15	98,50
IDV_te (3)	No detecta	5,50
IDV_te (4)	64	64,00
IDV_te (5)	11	29,38
IDV_te (6)	7	99,25
IDV_te (7)	1	100



IDV_te (8)	23	97,63
IDV_te (9)	No detecta	6,88
IDV_te (10)	79	49,63
IDV_te (11)	188	59,88
IDV_te (12)	7	99,00
IDV_te (13)	45	95,50
IDV_te (14)	1	100
IDV_te (15)	No detecta	9,88
IDV_te (16)	197	30,88
IDV_te (17)	27	83,63
IDV_te (18)	88	90,00
IDV_te (19)	No detecta	15,25
IDV_te (20)	87	46,00
IDV_te (21)	514	40,88
Media	<b>80,05</b>	<b>62,90</b>
Media sin fallos 3,9 y 15	<b>73,03</b>	<b>69,48</b>

Tabla 4: Datos obtenidos al aplicar PCA utilizando el estadístico  $T^2$

#### Resultados obtenidos utilizando el estadístico Q (SPE):

Tipo de fallo	Instante en el que se detecta el fallo	Alarmas detectadas (%)
IDV (1)	3	99,58
IDV (2)	7	98,75
IDV (3)	No detecta	20,20
IDV (4)	1	100
IDV (5)	3	62,50
IDV (7)	1	100
IDV (8)	11	97,70
IDV (9)	No detecta	17,71
IDV (10)	61	82,71
IDV (11)	30	79,17
IDV (12)	26	96,88
IDV (13)	12	98,54
IDV (14)	2	99,58
IDV (15)	No detecta	22,71
IDV (16)	19	64,38
IDV (17)	36	93,75
IDV (18)	77	86,88
IDV (19)	28	56,88



IDV (20)	34	74,79
IDV (21)	271	87,14
IDV_te (1)	161	100
IDV_te (2)	169	99,50
IDV_te (3)	633	25,75
IDV_te (4)	161	100
IDV_te (5)	161	49,75
IDV_te (6)	161	100
IDV_te (7)	161	100
IDV_te (8)	171	98,00
IDV_te (9)	No detecta	21,12
IDV_te (10)	195	76,12
IDV_te (11)	200	77,38
IDV_te (12)	169	98,63
IDV_te (13)	198	96,00
IDV_te (14)	162	99,25
IDV_te (15)	867	26,63
IDV_te (16)	176	71,13
IDV_te (17)	182	97,88
IDV_te (18)	236	93,38
IDV_te (19)	348	56,25
IDV_te (20)	241	73,63
IDV_te (21)	399	70,13
Media	<b>97,6</b>	<b>77,64</b>
Media sin fallos 3,9 y 15	<b>42,89</b>	<b>86,75</b>

Tabla 5: Datos obtenidos aplicando PCA con el estadístico Q (SPE)

En las tablas 4 y 5 se muestran los resultados obtenidos para cada tipo de fallo del sistema. Los fallos IDV corresponden con el bloque de datos de 480 observaciones y los fallos IDV\_te al bloque de datos de 960 mediciones, donde las primeras 160 se corresponden a un comportamiento sin fallo.

Se observa que con ambos estadísticos se pueden detectar la mayoría de los fallos, es decir que se supere el umbral 10 veces de forma consecutiva. Los fallos 3,9 y 15 denominados incipientes al ser muy difíciles de detectar, solamente son detectados con el estadístico Q y utilizando el bloque de datos de 960 observaciones. Además, con el estadístico Q se observa que se diagnostican un mayor número de fallos ya que el estadístico T no puede detectar el fallo 19 y los fallos 4 y 16 solo los detecta con los datos de 960 mediciones.



Si calculamos la media del porcentaje de alarmas detectadas por ambos estadísticos, el estadístico  $T^2$  detecta el 62,90 % y el estadístico Q 77,64%. Sin embargo, si calculamos la media del tiempo que tarda en detectarse un fallo, el estadístico  $T^2$  tarda de media un 80,05 y el estadístico Q 97,6. Por lo tanto, en general estadístico Q proporciona mayor seguridad para diagnosticar los fallos del proceso ya que detecta un mayor número de fallos y de anomalías producidas. Además, se han calculado las medias de los resultados obtenidos sin tener en cuenta los fallos 3,9 y 15. Se puede observar como las medias mejoran llegando hasta un 86,75% de anomalías detectadas mediante el estadístico Q y a un 69,48% con el estadístico  $T^2$ . El valor que más mejora es el tiempo detectado mediante el estadístico Q ya que disminuye de 97,6 a 42,89 debido a que los fallos eliminados se detectaban muy tarde.

#### 4.2.3. Contribución de las variables al fallo

En último lugar, se van a realizar los diagramas de contribución de cada variable a que se produzca el fallo en el sistema. De esta manera, podremos averiguar qué variables son las que más están afectando negativamente al proceso.

Para ello, como ya explicamos en el capítulo II la contribución se calcula diferente dependiendo el estadístico que estamos utilizando. Para el caso de Q, simplemente se calcula el vector residuo  $r$  (Ecuación 15) en el instante de tiempo de fallo y hacemos un diagrama de barras del cuadrado de ese vector para tener solo valores positivos. Por otro lado, para el estadístico  $T^2$  hay que calcular la matriz  $T$  (Ecuación 18), después hallar la contribución de cada variable a que aparezcan anomalías en las componentes principales (Ecuación 19) y finalmente sumar esta contribución (Ecuación 20).

Para este apartado, se van a utilizar dos fallos de ejemplo para visualizar los diagramas de barras de la contribución de las variables. Se utilizarán los



fallos 1 y 14 que tienen un alto porcentaje de alarmas detectadas por ambos estadísticos por lo que nos mostrarán datos más fiables a la hora de averiguar qué variables están contribuyendo más al desajuste del sistema.

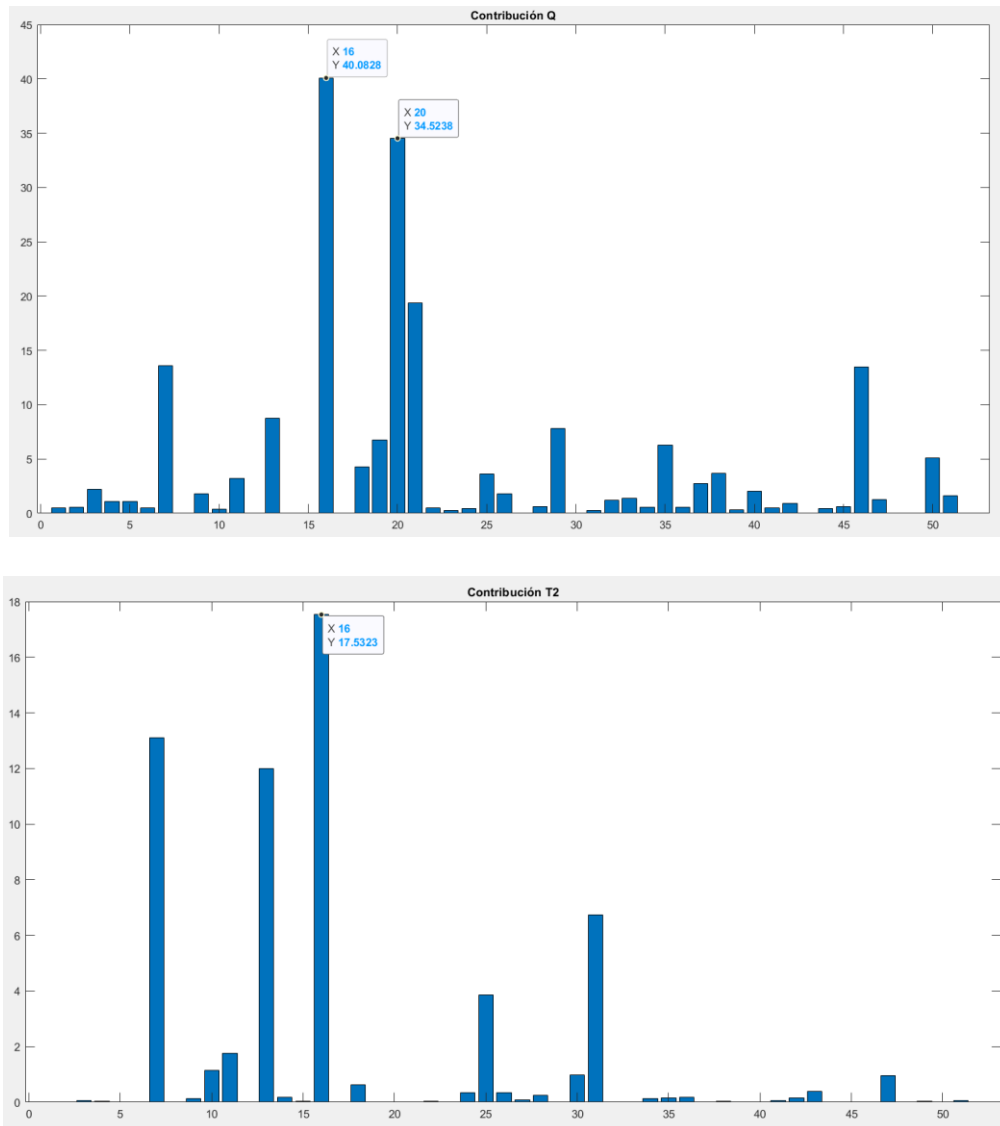


Figura 16: Contribución de cada variable a que se produzca el fallo 1

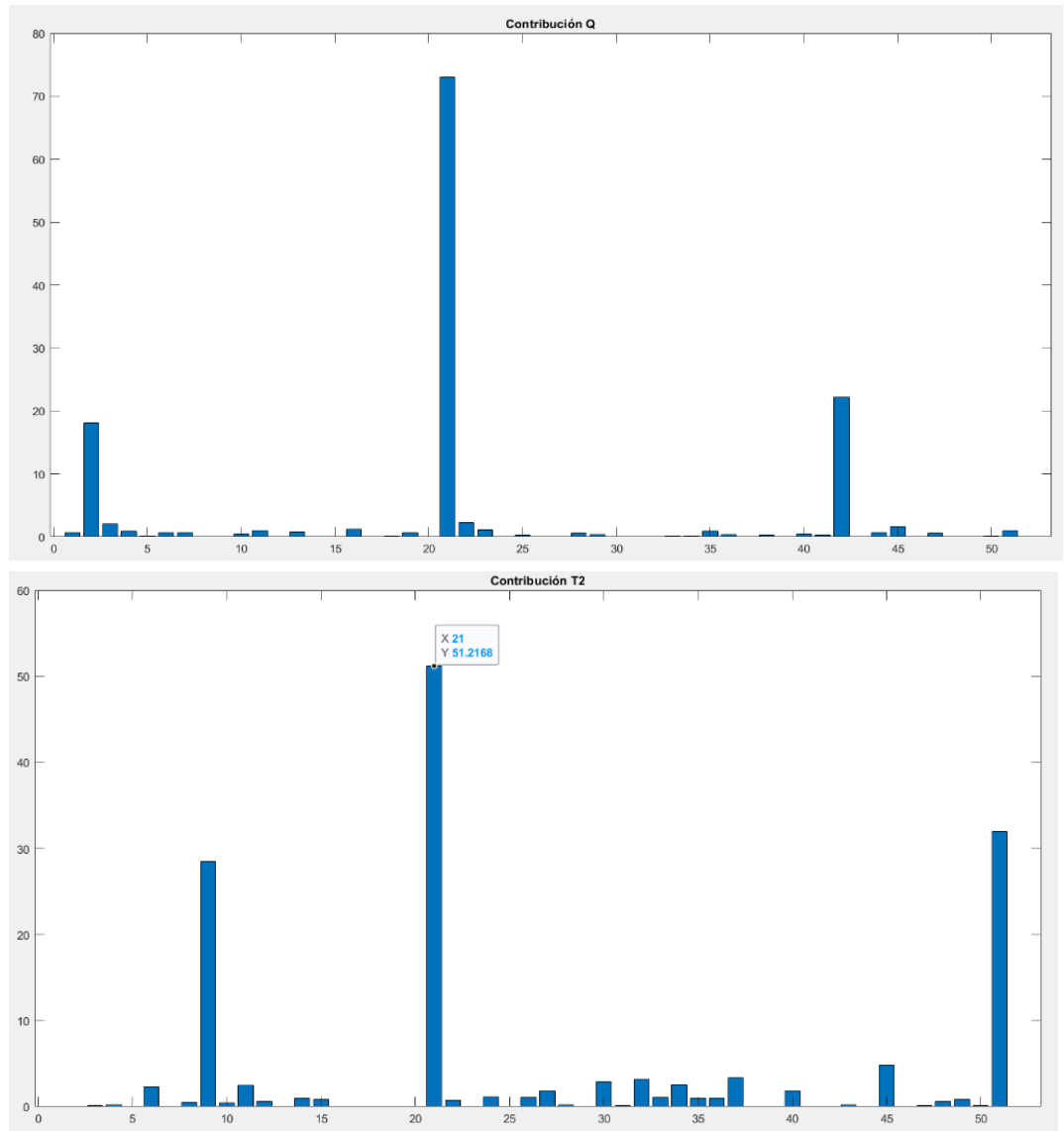


Figura 17: Contribución de cada variable a que se produzca el fallo 14





En la figura 16 se observa como las variables que más contribuyen al fallo 1 son la 16 (Presión del destilador) y la 20 (Potencia del compresor). La variable 16 es detectada a través de ambos estadísticos mientras que la 20 solo la detecta el estadístico Q. Por otro lado, para el fallo 14 observamos en la figura 17 que la variable que más contribuye al fallo es la 21 (Temperatura de la salida de agua de refrigeración del reactor) para ambos estadísticos.

Se ha realizado el estudio de los diagramas de contribución para todos los datos de simulación del proceso y así poder detectar que variables son las que suelen contribuir más a que se produzca un fallo en el sistema:

Tipo de fallo	Mayor contribución (T <sup>2</sup> )	Mayor contribución (Q)
IDV (1)	Variables <b>16</b>	Variables <b>16</b> y 20
IDV (2)	Variables <b>24 y 30</b>	Variables <b>16,24 y 30</b>
IDV (3)	No detecta	No detecta
IDV (4)	No detecta	Variables 9 y 51
IDV (5)	Variables 7,11,13,16	Variables 29
IDV (7)	Variables 7,13 y <b>16</b>	Variables <b>16</b>
IDV (8)	Variables <b>16</b>	Variables <b>16</b>
IDV (9)	No detecta	No detecta
IDV (10)	Variables <b>18</b>	Variables <b>18</b>
IDV (11)	Variables <b>51</b>	Variables 9 y <b>51</b>
IDV (12)	Variables 7, <b>11</b> y 16	Variables <b>11</b> y 19
IDV (13)	Variables 7,13, <b>16</b>	Variables <b>16</b>
IDV (14)	Variables <b>21</b>	Variables <b>21</b>
IDV (15)	No detecta	No detecta
IDV (16)	No detecta	Variables 50
IDV (17)	Variable <b>21</b>	Variable <b>21</b>
IDV (18)	Variables 7, <b>11</b> ,13 y 16	Variables <b>11</b> y 22
IDV (19)	No detecta	Variables 22,27 y 46
IDV (20)	Variable <b>46</b>	Variable <b>46</b>



IDV (21)	Variables 19 y 50	Variables 25 y 36
IDV_te (1)	Variables <b>16</b>	Variable <b>16</b> y 20
IDV_te (2)	Variables 10, <b>30</b> ,47	Variables <b>30</b>
IDV_te (3)	No detecta	Variable 25
IDV_te (4)	Variables 7,11,13	Variable 9,51
IDV_te (5)	Variables 13 y 16	Variables 11,18 y 30
IDV_te (6)	Variable 20	Variables 1 y 44
IDV_te (7)	Variable <b>16</b>	Variable <b>16</b>
IDV_te (8)	Variables <b>16</b>	Variables <b>16</b>
IDV_te (9)	No detecta	No detecta
IDV_te (10)	Variable 7,13, <b>16 y 19</b>	Variable <b>16 y 19</b>
IDV_te (11)	Variables <b>9 y 51</b>	Variable <b>9 y 51</b>
IDV_te (12)	Variable 14, <b>25</b> y 43	Variable <b>25</b>
IDV_te (13)	Variable <b>16</b>	Variable <b>16</b>
IDV_te (14)	Variable 9 y 51	Variable 21
IDV_te (15)	No detecta	Variable 20 y 21
IDV_te (16)	Variable <b>18</b> y 19	Variable <b>18,25</b> y 50
IDV_te (17)	Variable <b>21</b>	Variable <b>21</b>
IDV_te (18)	Variable <b>22</b>	Variables <b>11 y 22</b>
IDV_te (19)	No detecta	Variable 20
IDV_te (20)	Variable 13 y <b>46</b>	Variables <b>46</b>
IDV_te (21)	Variables 19 y 50	Variable 27

Tabla 6: Resultados de las variables que más contribuyen a que produzca fallo en el sistema

Se pueden observar en la tabla 6 los resultados tras haber realizado los cálculos mencionados para obtener las variables que más contribuyen a que se produzca fallo en el proceso. Se han resaltado en negrita las variables que se repiten tanto para el estadístico Q como el estadístico T2. Para algunos fallos no es posible detectar que variable es la que más ha contribuido ya que



el fallo no se detecta por lo que no podemos saber el tiempo de fallo necesario para el cálculo. En varios casos las variables con mayor contribución son diferentes dependiendo el estadístico que estamos utilizando, sin embargo, se puede observar que en la mayoría de las ocasiones se repiten no solo para ambos estadísticos sino entre los distintos bloques de datos utilizados. De esta manera, podemos comprobar que es un método bastante estable para poder clasificar y detectar cual son las variables que están perjudicando al sistema.

Para este caso, se ha observado que las variables que más se repiten son la 16 (Flujo de alimentación de A y C), 21 (Temperatura del reactor) y 46 (Concentración de F en la Purga). Gracias a este método se podría revisar por qué estas variables afectan de forma negativa al sistema y si es necesario realizar ciertos ajustes en el proceso para que no se produzcan las anomalías que hemos detectado.

### **4.3. Uso del método t-SNE para la detección de fallos producidos en el sistema**

El segundo método que vamos a utilizar para detectar y diagnosticar las anomalías del proceso es la técnica no lineal de incrustación de vecinos estocásticos distribuidos en  $t$ , denominado de manera más breve por las siglas t-SNE (*t-distributed stochastic neighbor embedding*). Para el empleo de esta técnica se va a utilizar el bloque de datos de 960 observaciones ya que en el método anterior comprobamos que los resultados para ambos bloques de datos eran muy similares, por tanto, se ha elegido el bloque que más datos ofrece y que además contiene tanto datos de comportamiento normal (los 160 primeros) como datos con fallo del sistema (los 800 restantes).



#### 4.3.1. Detección de anomalías del proceso usando la distancia euclídea

En primer lugar, se crea el modelo t-SNE de los datos de comportamiento normal. En este caso la matriz de datos originales será de dimensión  $960 \times 52$ . Los datos se normalizan a media cero y varianza uno para obtener unos mejores resultados como hicimos con PCA anteriormente.

En la técnica t-SNE se pueden utilizar distintos tipos de distancias para medir la separación entre puntos pertenecientes al espacio de estudio. Matlab por defecto utiliza la distancia euclidiana, sin embargo, también nos permite trabajar con varios tipos de distancias diferentes ya sea la distancia chebychev, minkowski, mahalanobis, cosine, etc. Para la primera prueba con este método vamos a utilizar la distancia euclidiana y se va a utilizar la función *tsne* que viene incorporada directamente en Matlab y nos calcula la matriz  $Y$  de dimensión reducida. Por defecto, esta función además de trabajar con la distancia euclidiana reduce el espacio original a un espacio de dos dimensiones pudiendo ser ajustado para que se reduzca a un espacio de tres dimensiones. De momento se utilizarán los valores por defecto de la función y se observarán los resultados obtenidos. De esta manera, se convertirá la matriz original  $X \in \mathfrak{R}^{960 \times 52}$  a una matriz  $Y \in \mathfrak{R}^{960 \times 2}$  donde  $y_i$  serán los datos proyectados de  $x_i$  con los que trabajaremos a partir de ahora.

Teniendo ya la nueva matriz  $Y$  reducida, se procede a calcular el parámetro  $A \in \mathfrak{R}^{960 \times 2}$  (Ecuación 25) que nos permite relacionar el espacio de datos de alta dimensión con el nuevo espacio de baja dimensión. Además, también se calculará el parámetro  $R = \left( \frac{Y \cdot Y^T}{n-1} \right)$  (ecuación 26) que nos ayudará a calcular el estadístico  $T^2$  para analizar si los datos del sistema presentan alguna anomalía.

Una vez se obtienen estos parámetros, se calcula el estadístico  $T^2$  (Ecuación 29) y se establece el umbral observando el gráfico de los datos del estadístico

calculados para los datos de comportamiento normal. En este caso se ha decidido utilizar un umbral  $T_a^2 = 6$ , se puede observar en la figura 18 que prácticamente todos los valores del estadístico  $T^2$  se encuentran por debajo de este valor, que es lo que debería ocurrir al trabajar con los datos sin fallo en el sistema.

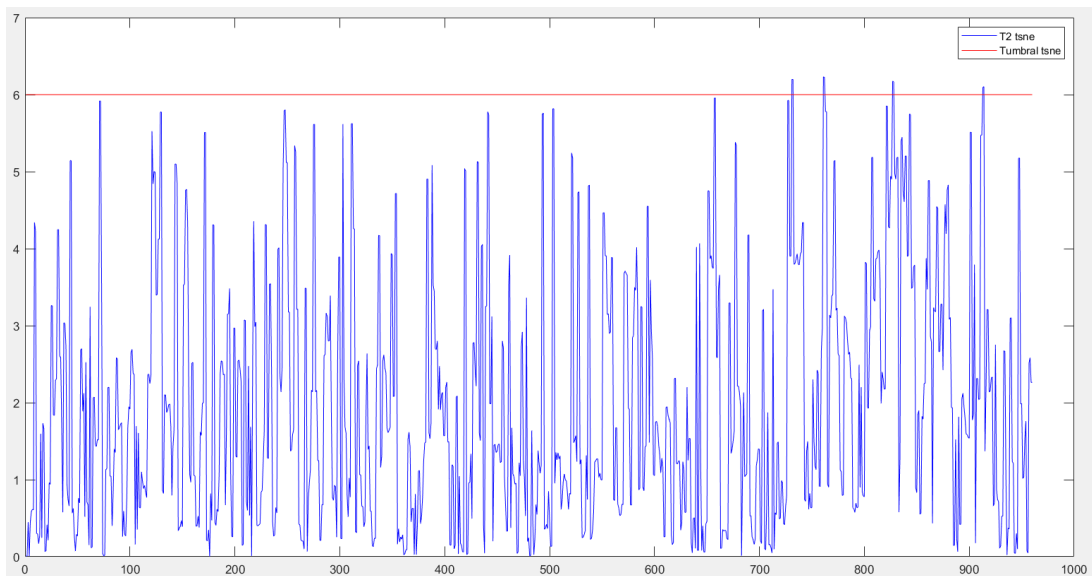


Figura 18: Comparación de los valores calculados del estadístico  $T^2$  con el umbral  $T_a^2$  para un comportamiento normal del proceso

Después de haber realizado todos los cálculos necesarios para los datos de comportamiento normal, se procede a utilizar las simulaciones que contienen los datos con fallos del proceso para comprobar si el método detecta las anomalías que se producen.

Para ello, primero se importa la matriz principal  $X \in \mathbb{R}^{960 \times 52}$  compuesta por las 52 variables medidas y las 960 observaciones, donde el fallo se produce a partir de la número 161. Al igual que en el paso anterior, se normaliza esta matriz a media cero y varianza uno. Para calcular la matriz  $Y$  de dimensión reducida se va a utilizar la matriz  $A$  calculada anteriormente para obtener la matriz  $Y$  a través de la ecuación 26.

Una vez obtenida la matriz  $Y$  de dimensión reducida, se procede a calcular el estadístico  $T^2$  a través de la ecuación 29 utilizando el parámetro  $R$  calculado también anteriormente y la nueva matriz  $Y$  que acabamos de obtener.

Por último, nos queda comparar los valores obtenidos con el umbral  $T_a^2$  ya fijado. Se utilizará la misma metodología que con el método PCA para analizar estos datos, si los datos del estadístico  $T^2$  superan el umbral seleccionado 10 veces consecutivas, se considera que ha habido fallo del sistema. Además, se hará también un recuento del número de alarmas detectadas, es decir, el número total de veces que se detecta que ha habido una anomalía en el sistema porque se haya superado el umbral  $T_a^2$ .

Al igual que con PCA, para mostrar una prueba del gráfico del estadístico  $T^2$  del sistema se van a elegir los fallos 1,8 y 15 para comprobar que tienen una forma similar a las figuras 13,14 y 15 calculadas mediante PCA respectivas a estos fallos.

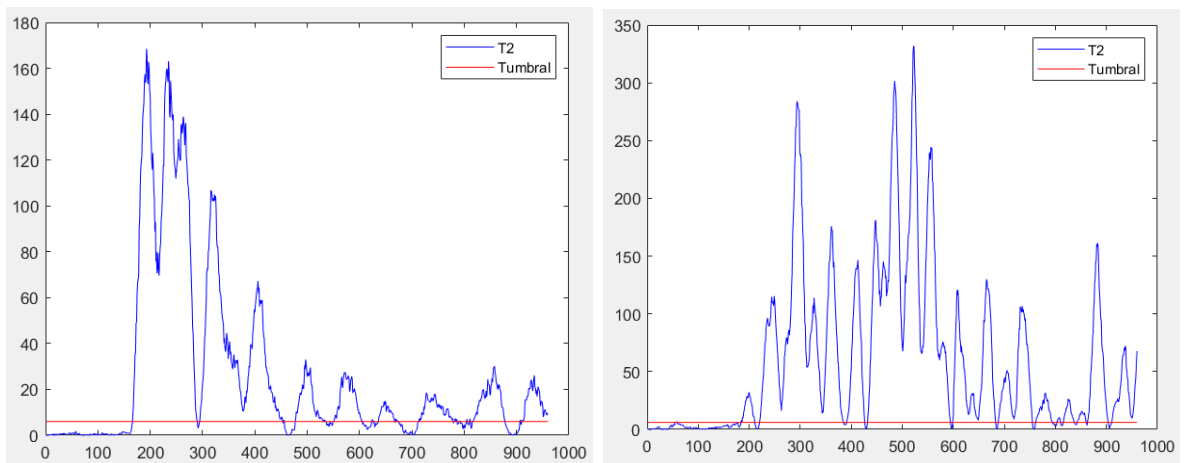


Figura 19: Representación del estadístico  $T^2$  para los fallos 1 (izquierda) y 8 (derecha) del sistema mediante t-SNE

En la figura 19 se puede observar gráficamente los valores del estadístico  $T^2$  de los fallos 1 y 8 del sistema. También se ha representado el umbral  $T_a^2$  para lograr una fácil interpretación visual de los resultados obtenidos. Si se comparan estos resultados con los mostrados en las figuras 13 y 14 para los



fallos 1 y 8 calculados con PCA, se comprueba que los valores del estadístico tienen la misma evolución a lo largo de la gráfica en ambos métodos.

Para el fallo 1, los valores superan rápidamente el umbral  $T_{\alpha}^2$  una vez se supera la observación 160 donde comienza a fallar el sistema llegando a valores muy altos del estadístico que se va reduciendo y estabilizando a lo largo del tiempo aunque al contrario que PCA, al estabilizarse hay tramos donde deja de detectar anomalías en el sistema, lo cual significa un menor número de alarmas detectadas.

Por otra parte, en el gráfico de la derecha para el fallo 8 también se observa una evolución similar al mostrado en la figura 14 con el método PCA. El primer tramo de comportamiento normal sin fallo se ve reflejado en valores del estadístico  $T^2$  menores que el umbral en todo momento, cuando el sistema comienza a fallar se tarda poco en obtener valores altos de  $T^2$  que superan en gran medida el umbral. Al igual que sucedía con PCA, estos valores no se estabilizan como en el fallo 1 sino que se observan subidas y bajadas continuas a lo largo del tiempo, estando la mayoría por encima del umbral fijado.

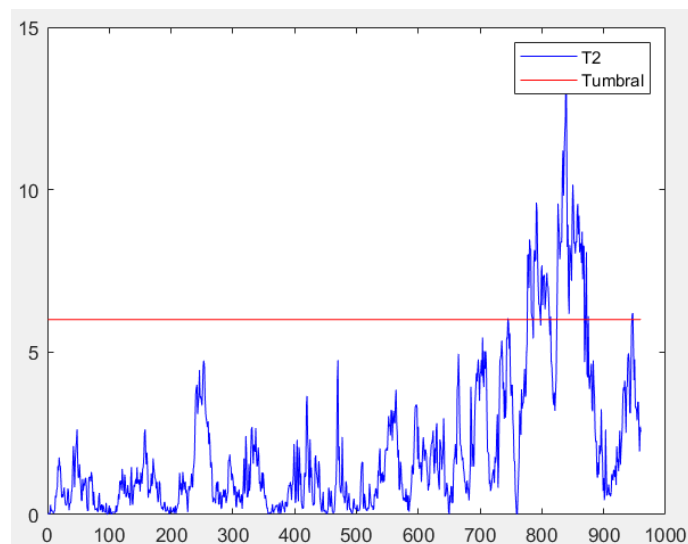


Figura 20: Representación del estadístico  $T^2$  para el fallo 15 del sistema mediante t-SNE



También se ha representado el fallo 15 del sistema que como ya se ha comentado, es un fallo que posee unas características más especiales al ser muy difícil de detectar por técnicas de control estadístico de procesos. Mediante PCA no era posible detectar el fallo ya que no se superaba en ningún momento el umbral 10 veces de forma consecutiva. Sin embargo, a través de t-SNE sí que ha sido posible detectar el fallo aunque mucho más tarde del momento en el que comienza a producirse que es a partir de la observación 160. Además, se logran muy pocas alarmas detectadas a pesar de detectar el fallo. Por tanto, este método tampoco es realmente eficaz para detectar estos fallos tan complejos de la planta.

Finalmente, se ha recopilado todos los datos que se han obtenido al analizar todos las simulaciones del proceso ante todos los posibles fallos que pueden ocurrir en la planta Tennessee Eastman.

Tipo de fallo	Tiempo de fallo	Alarmas detectadas (%)
d01_te	5	84,75
d02_te	25	97,75
d03_te	7	27,875
d04_te	59	13,875
d05_te	15	97,125
d06_te	1	100
d07_te	1	55
d08_te	7	94,375
d09_te	5	25
d10_te	9	63,5
d11_te	174	24,375
d12_te	18	92,5
d13_te	45	93,25
d14_te	No detecta	14,25
d15_te	615	23,5
d16_te	35	57,25
d17_te	41	56,25
d18_te	99	90,25





d19_te	68	10,25
d20_te	73	58,125
d21_te	577	22,375
Media	<b>93,95</b>	<b>57,22</b>
Media sin fallos 3,9 y 15	<b>73,64</b>	<b>62,51</b>

Tabla 7: Datos obtenidos aplicando t-SNE con los valores por efecto en Matlab

En primer lugar, lo que más destaca es que es la capacidad del método por detectar los fallos 3 y 9 (Temperatura de alimentación D) y 15 (Válvula del agua de refrigerante del condensador) que son los más complejos para el control estadístico de esta planta. Sin embargo, se observa que el porcentaje de alarmas detectadas para estos fallos está en torno a un 25%, un valor pequeño ya que no se detectan tres cuartas partes de las anomalías que se producen en el sistema en estos casos. Además, también destaca que el único fallo que no se llega a detectar es el fallo 14 (Bloqueo de la válvula del agua de refrigerante del reactor) con un porcentaje muy pequeño de anomalías detectadas (14,25 %) mientras que PCA detectaba este fallo muy temprano y prácticamente todas anomalías que ocasionaba al sistema.

Para medir el tiempo que se tarda en detectar el fallo se ha contabilizado a partir de la observación 160 que es cuando empieza a fallar el sistema, es decir si se detecta en el instante 1 significa que se está detectando en la observación 161. La media de tiempo que tarda en detectar los fallos el método t-SNE es 93,95 un valor muy parecido a los proporcionados por el método PCA. Sin embargo, el porcentaje de alarmas detectadas en el sistema es 57,22% un valor menor que cuando utilizamos PCA. Es decir, ambos métodos funcionan de forma similar pero PCA consigue detectar mayores anomalías en el proceso. Además, se ha hecho la media sin los fallos 3,9 y 15 y se obtienen mejores resultados debido a que estos errores son difíciles de detectar.



### 4.3.2. Detección de anomalías del proceso mediante reducción a 3 dimensiones con la técnica t-SNE

Hemos comprobado que para los valores por defecto que proporciona Matlab para la función *tsne* se obtienen resultados algo peores que utilizando el método PCA ya que se detectan menos porcentaje de anomalías en el proceso.

Vamos a utilizar la misma metodología pero cambiando los diferentes parámetros que tiene la formula *tsne* para intentar lograr mejores resultados. En primer lugar, vamos a seguir utilizando la misma distancia entre puntos por defecto pero en esta ocasión, se va a reducir la dimensionalidad de la matriz  $X$  inicial  $X \in \mathbb{R}^{960 \times 52}$  a una matriz  $Y \in \mathbb{R}^{960 \times 3}$ , es decir, se va a reducir las 52 variables iniciales a una dimensión de 3 ya que mediante la técnica t-SNE podemos elegir si reducir a 2 o 3 dimensiones, siendo 2 la que utiliza Matlab por defecto. El inconveniente es que se va a trabajar con un mayor número de datos al haber una columna a mayores con 960 valores más, así comprobaremos si este aumento de datos mejora la capacidad del método para detectar los fallos producidos en la planta.

En lo referente a como aplicamos la técnica t-SNE, el proceso es el mismo que en el caso anterior. Se utilizan en primer lugar los datos de comportamiento normal sin fallo del proceso normalizados a media cero y varianza uno y se calcula la matriz  $Y$  de dimensión reducida, donde en este caso se impondrá la condición de contenga tres dimensiones. Una vez Matlab nos calcula la matriz  $Y$ , simplemente tenemos que calcular los parámetros  $A$  (Ecuación 25) y  $R$  (Ecuación 26) que hemos utilizado en el apartado anterior. Lo último que vamos a comprobar es si el umbral  $T_a^2$  anterior establecido a un valor de 6, sigue siendo válido en este caso. Calculamos el estadístico para todas las 960 observaciones de la matriz  $Y$ , en la figura 21 se observa como la mayoría de los datos siguen por debajo del umbral y por tanto se puede seguir utilizando este valor de referencia.

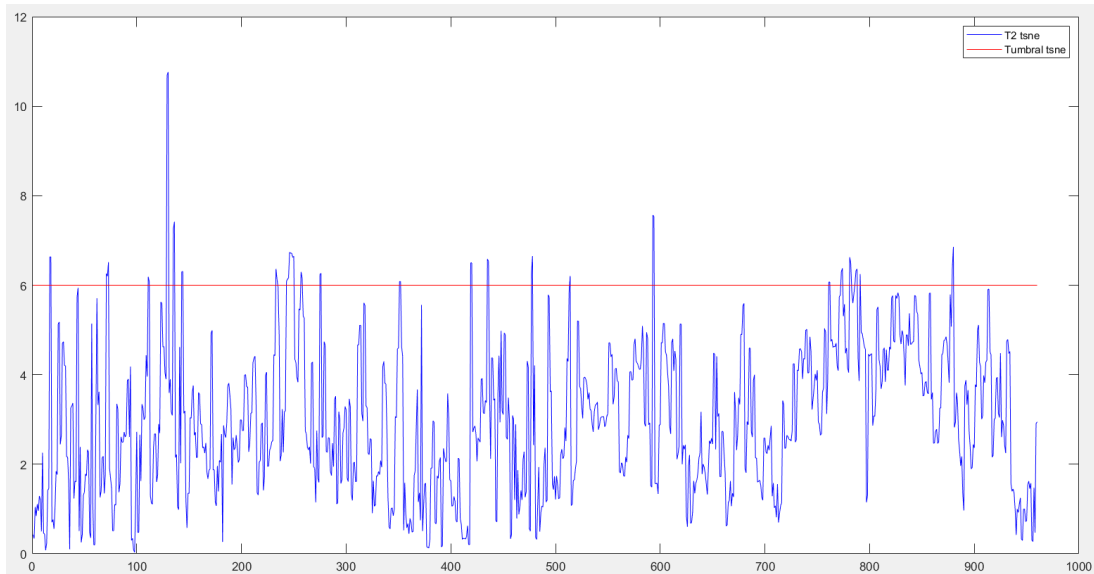


Figura 21: Comparación de los valores calculados del estadístico  $T^2$  para una matriz  $Y$  reducida a 3 dimensiones con el umbral  $T_\alpha^2$

Después de obtener los parámetros  $A$ ,  $R$  y el umbral del estadístico  $T^2$ , se procede a calcular las matrices  $Y$  reducidas a 3 dimensiones de la misma manera que antes para los datos con fallo de la planta química. La única novedad es que la matriz  $A$  que hemos calculado en esta ocasión tiene dimensión  $52 \times 3$  para que las nuevas matrices  $Y$  tengan la dimensión deseada al calcularse a través de la ecuación 26. Finalmente, se utilizan estas nuevas matrices junto al parámetro  $R$  de los datos de comportamiento normal para calcular el estadístico  $T^2$  y compararlo con el umbral  $T_\alpha^2$  para detectar las anomalías del sistema. Como hemos hecho en todos los casos, podremos afirmar que ha habido fallo si se supera el umbral  $T_\alpha^2$  10 veces consecutivas y cada vez que se supere el umbral se contará como alarma detectada. Los datos que se han obtenido han sido los siguientes:



Tipo de fallo	Tiempo de fallo	Alarmas detectadas (%)
IDV_te (1)	6	99,375
IDV_te (2)	25	77,875
IDV_te (3)	42	17,125
IDV_te (4)	151	7,875
IDV_te (5)	14	98,25
IDV_te (6)	1	100
IDV_te (7)	2	52
IDV_te (8)	7	95
IDV_te (9)	5	14,125
IDV_te (10)	13	60,625
IDV_te (11)	177	19,625
IDV_te (12)	15	95,75
IDV_te (13)	47	94,25
IDV_te (14)	No detecta	5
IDV_te (15)	615	19
IDV_te (16)	191	49,25
IDV_te (17)	57	39,5
IDV_te (18)	94	89,25
IDV_te (19)	No detecta	3,875
IDV_te (20)	79	57,875
IDV_te (21)	641	21,625
Media	<b>114,84</b>	<b>53,20</b>
Media sin los fallos 3, 9 y 15	<b>95,00</b>	<b>59,28</b>

Tabla 8: Datos obtenidos aplicando t-SNE reduciendo a 3 dimensiones

Después de realizar la recopilación de los resultados para todos los fallos del sistema, se observa que al igual que en el caso anterior los fallos 3,9 y 15 sí que se detectan pero con un porcentaje de alarmas detectadas menor que anteriormente. Además, el fallo 14 se sigue sin detectar y en esta ocasión tampoco se detecta el fallo 19.

Para el tiempo de fallo de la tabla 8 también se ha puesto el tiempo a partir del instante 160. Si realizamos la media del tiempo que se tarda en detectar fallo nos sale 114,84 para este caso cuando en el apartado anterior el valor era de 93,95. Por otro lado, el porcentaje alarmas detectadas es 53,20 y antes habíamos obtenido un porcentaje de anomalías detectadas de 57,22.



Comparando los resultados se observa que no hemos mejorado la detección de anomalías y fallos del sistema en comparación con el caso anterior por lo que probaremos a realizar otros ajustes diferentes para mejorar estos resultados. Al igual que en el caso anterior, se ha calculado la media sin los fallos 3,9 y 15, sin embargo, aunque se obtienen valores algo más elevados siguen siendo peores que anteriormente.

#### **4.3.3. Detección de anomalías del proceso con la técnica t-SNE utilizando distintas distancias entre puntos**

En el apartado anterior probamos a cambiar el número de dimensiones que Matlab reducía a 2 por defecto al utilizar la función *tsne*, sin embargo, no se obtuvieron mejores resultados que para el caso de todas las opciones por defecto ni para los datos que se han obtenido utilizando el método PCA.

Por tanto, en esta ocasión se va a probar a cambiar la distancia entre puntos que utiliza Matlab de manera predeterminada. Para ello, primero se va a utilizar la función *gscatter* que nos ayudará a visualizar cual pueden ser las mejores distancias para aplicar el método t-SNE al sistema que estamos utilizando. Para utilizar esta función se van a juntar todas las matrices  $\in \mathbb{R}^{480 \times 52}$  de datos de fallo en una sola matriz  $X \in \mathbb{R}^{8160 \times 52}$ , es decir, una matriz que contenga todos los datos con anomalías en el sistema para las 52 variables. Se ha utilizado el bloque de datos de 480 observaciones porque queremos tener solo los datos con fallo y en el caso del bloque de simulaciones con 960 las 160 primeras son de comportamiento normal de la planta y además al ser muchos más datos sería más difícil interpretar el gráfico que vamos a obtener mediante la función *gscatter*. Además, no se han utilizado los fallos 3,9 y 15 por ser fallos que al ser más costosos de detectar, dificultarían la visualización del gráfico obtenido. Tampoco se ha utilizado el fallo 6 porque no hay datos de este fallo para este bloque de simulaciones.

Para utilizar la función `gscatter` simplemente hay que aplicar la función `tsne` a la matriz  $X \in \mathbb{R}^{8160 \times 52}$  para transformarla a una matriz  $Y \in \mathbb{R}^{8160 \times 2}$  con dos dimensiones que será la que utilizaremos para el gráfico. Finalmente se crea una matriz que denominaremos *clases* que asigne un número similar a los datos que provengan de la misma simulación, permitiéndonos que en el gráfico cada bloque de datos tenga un color diferente y así poder observar si se separa correctamente. Los gráficos obtenidos son los siguientes:

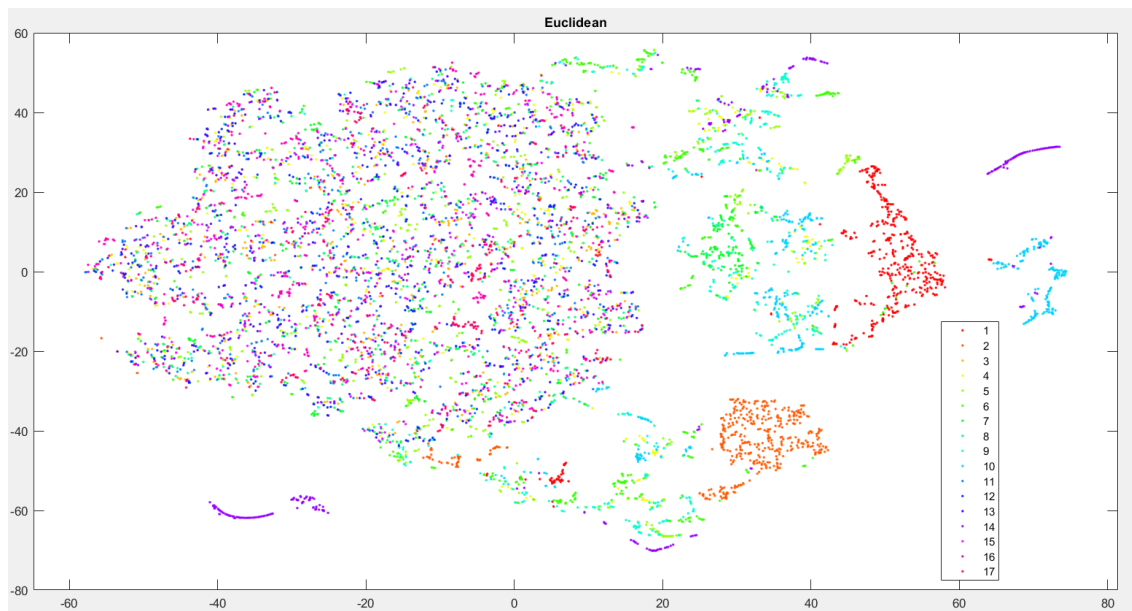


Figura 22: Gráfico obtenido mediante la función `gscatter` y la distancia euclídea

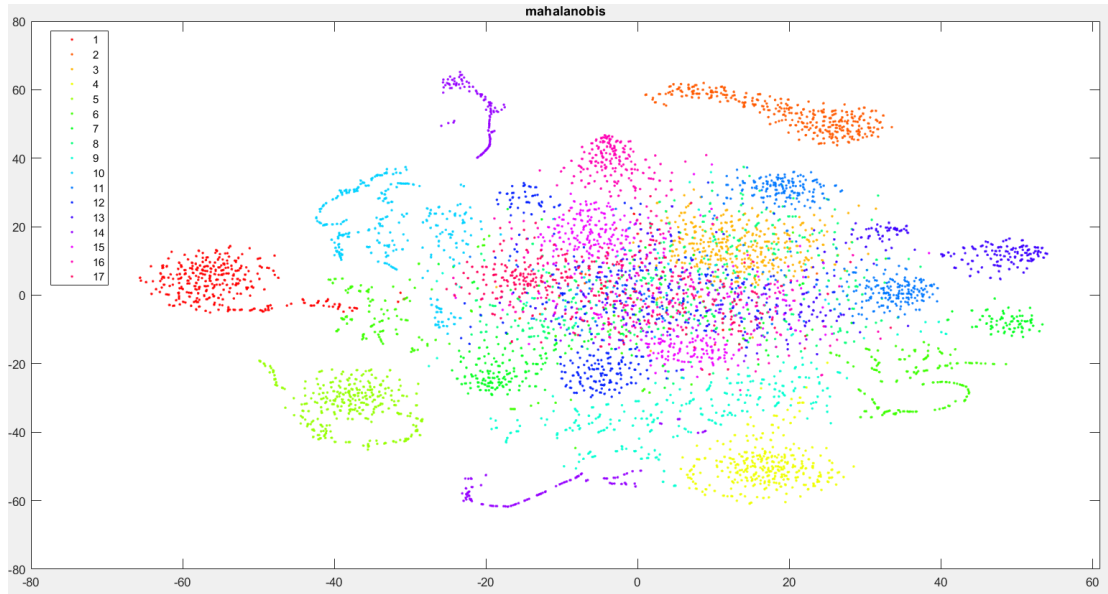


Figura 23: Gráfico obtenido mediante la función gscatter y la distancia mahalanobis

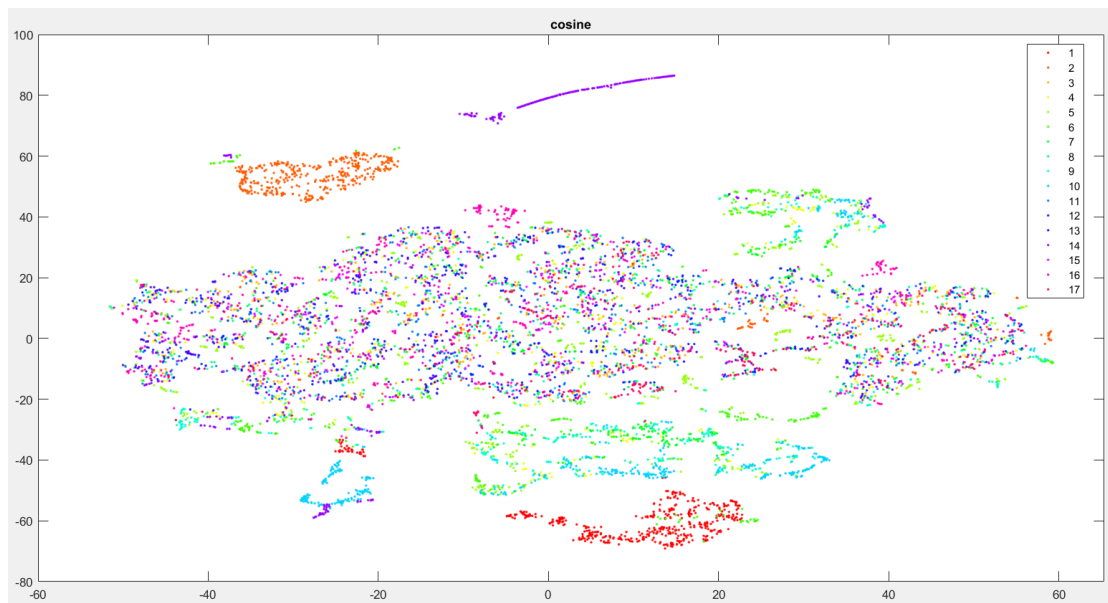


Figura 24: Gráfico obtenido mediante la función gscatter y la distancia coseno

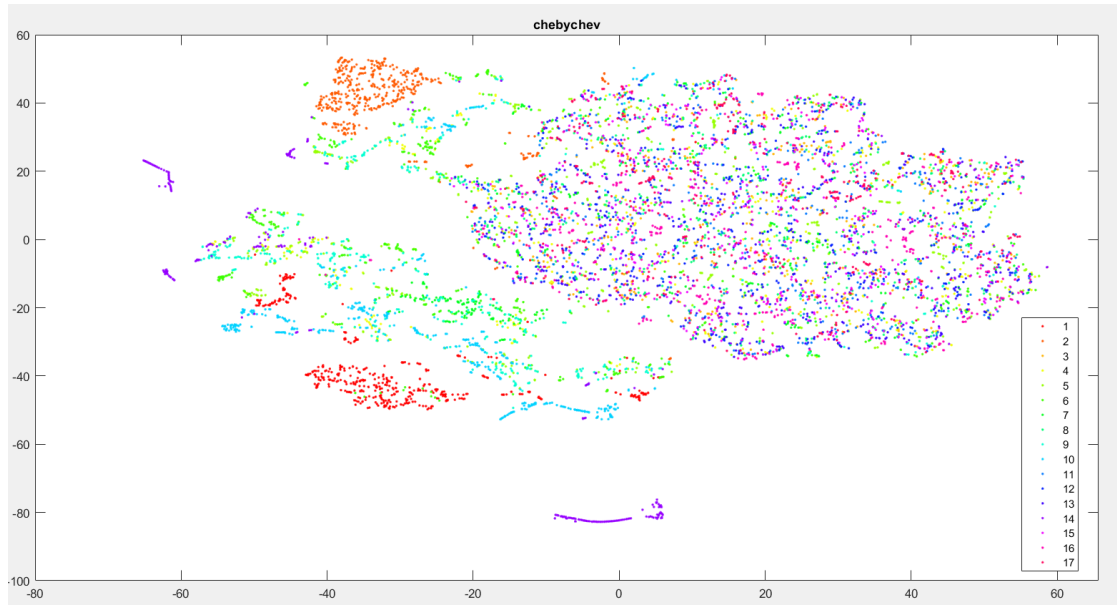


Figura 25: Gráfico obtenido mediante la función gscatter y la distancia chebychev

Vamos a comenzar por descartar cual serían las peores distancias para detectar fallos de la planta, analizando los gráficos correspondientes a las figuras 22, 23, 24 y 25. En todas las imágenes se observan que hay varios grupos que se diferencian sobre el resto como son el grupo 1, 2, 9 o 17 entre otros, que se posicionan alejados y agrupados en una parte del gráfico, lo que significa que la mayoría de distancias puede identificar bien esos grupos de datos correspondientes a un fallo de la planta en concreto. Sin embargo, observamos que la distancia euclidiana a simple vista es la que peor parece que separa estos datos, siendo esta la distancia que hemos utilizado por defecto para conseguir los resultados anteriores. A continuación, le siguen las distancias cosine y chebychev que separan algún grupo más de la mayoría de datos pero se asemejan al gráfico de la distancia euclidiana. Por otra parte, se encuentra el gráfico correspondiente a la distancia mahalanobis en la Figura 24, el cual destaca positivamente frente a los demás. Esto es debido a que visualmente observamos como utilizando esta distancia se separan la mayoría de los grupos perteneciente a cada fallo del sistema.





A priori parece que la distancia mahalanobis sería la más eficaz para detectar las anomalías de la planta que se está estudiando. Se va a comprobar si los datos que nos están proporcionando estos gráficos son correctos. Para ello, vamos a aplicar la técnica t-SNE con todas estas distancias salvo la euclidiana que ya la hemos analizado. Se va a reducir a la dimensionalidad por defecto establecida en Matlab que será a 2 dimensiones.

Aplicando el mismo procedimiento que en los casos anteriores, se obtienen los siguientes datos:

Tipo de fallo	Distancia Mahalanobis		Distancia Coseno		Distancia Chebychev	
	Tiempo de fallo	Alarmas detectadas (%)	Tiempo de fallo	Alarmas detectadas (%)	Tiempo de fallo	Alarmas detectadas (%)
IDV_te (1)	3	95	4	99,625	8	42,75
IDV_te (2)	23	97,625	26	94,125	15	28,125
IDV_te (3)	No detecta	3,5	411	6,625	No detecta	4,5
IDV_te (4)	No detecta	3,75	75	6	No detecta	3
IDV_te (5)	1	100	27	96,75	1	98,875
IDV_te (6)	1	100	1	100	17	98,5
IDV_te (7)	1	35,75	2	40,125	2	53,625
IDV_te (8)	22	86,875	7	91,875	31	82,875
IDV_te (9)	No detecta	3,625	225	8	No detecta	3,125
IDV_te (10)	28	63,25	28	40,75	27	32,125
IDV_te (11)	293	7,625	184	16,75	303	6,875
IDV_te (12)	11	95,5	65	86,125	22	90,75
IDV_te (13)	38	88,5	45	95	78	83,875
IDV_te (14)	No detecta	4,5	No detecta	5,75	713	57,125
IDV_te (15)	No detecta	10,625	627	17,25	No detecta	7



IDV_te (16)	35	54,5	265	35,5	245	30,125
IDV_te (17)	80	39,5	49	34,5	27	64,75
IDV_te (18)	90	90	111	85,75	101	88,125
IDV_te (19)	No detecta	6,875	No detecta	5,125	No detecta	0,875
IDV_te (20)	67	82,25	85	27	79	37,625
IDV_te (21)	491	46,375	464	45,875	723	6,75
Media	<b>78,93</b>	<b>53,13</b>	<b>142,16</b>	<b>49,45</b>	<b>149,50</b>	<b>43,875</b>
Media sin los fallos 3,9 y 15	<b>78,93</b>	<b>61,02</b>	<b>89,88</b>	<b>55,92</b>	<b>149,50</b>	<b>50,38</b>

Tabla 9: Resultados obtenidos aplicando t-SNE utilizando distintas distancias

A la vista de los resultados obtenidos se pueden destacar varios datos interesantes. Las distancias Coseno y Chebychev aportan los peores datos obtenidos hasta el momento de todos los métodos utilizados, siendo Chebychev la que menor porcentaje de alarmas detecta y la que más tarde detecta los fallos de la planta. Por otro lado, habíamos supuesto a priori que la distancia Mahalanobis iba a ser la que mejor resultados nos iba a otorgar, lo cual se cumple pero con algún inconveniente ya que no se detectan ni los fallos 3 y 9 en ningún momento al contrario que pasaba con la distancia euclidiana que si los detectaba aunque con poco porcentaje de alarmas detectadas. Además, tampoco detecta ni los fallos 4, 14 y 19 los cuales tampoco se detectan demasiado bien con ninguna de las distancias. Un dato positivo que destacar es que la distancia Mahalanobis detecta bastante rápido los fallos en comparación a las otras distancias utilizadas. Si se comparan las medias sin los fallos 3,9 y 15 el resultado es similar a los anteriores, mejorando los valores y manteniéndose la distancia mahalanobis como la distancia que mejores resultados proporciona.



#### 4.3.4. Mejora de la detección de anomalías del proceso utilizando la distancia Mahalanobis

Ya hemos analizado diversas distancias y la que mejor resultados nos ha otorgado para detectar las anomalías del proceso ha sido la distancia Mahalanobis. A través del gráfico de la figura 24 ya pudimos hacernos a la idea de que era una distancia que iba a ser capaz de detectar la mayoría de los fallos de la planta de forma correcta. Después se realizaron los cálculos y se confirmó que era la distancia que mejores resultados nos iba a proporcionar aunque le sigue muy de cerca la distancia euclidiana que utiliza Matlab por defecto.

Vamos a tratar de mejorar los resultados obtenidos con la técnica mahalanobis. Para ello, vamos a reducir el número de variables de la planta que estamos utilizando de 52 a 33 variables. Se va a utilizar por tanto una matriz inicial X que contenga las variables de la 1 a la 22 y desde la 42 hasta la 52. Es decir, se eliminan desde la variable 23 hasta la variable 41. Se va a utilizar la misma metodología de siempre para obtener el tiempo que se tarda en detectar el fallo y el porcentaje de anomalías que se logran detectar.

Los resultados obtenidos son los siguientes:

Tipo de fallo	Tiempo de fallo	Alarmas detectadas (%)
IDV (1)	2	98,13
IDV (2)	37	98,00
IDV (3)	No detecta	3,75
IDV (4)	358	10,13
IDV (5)	9	99,00
IDV (6)	1	100
IDV (7)	2	44,63
IDV (8)	20	88,13
IDV (9)	No detecta	2,25



IDV (10)	100	35,00
IDV (11)	385	17,13
IDV (12)	16	91,88
IDV (13)	43	92,13
IDV (14)	21	89,5
IDV (15)	No detecta	5,5
IDV (16)	193	26,38
IDV (17)	27	74,25
IDV (18)	95	88,25
IDV (19)	No detecta	36,5
IDV (20)	97	47,5
IDV (21)	410	55,13
Media	<b>106,82</b>	<b>57,29</b>
Media sin los fallos 3,9 y 15	<b>106,82</b>	<b>66,20</b>

Tabla 10: Resultados obtenidos aplicando t-SNE utilizando 33 variables y la distancia mahalanobis

En primer lugar, se observa que se siguen sin detectar los fallos 3,9 y 15 pero no supone demasiado problema ya que no nos queremos centrar en detectar estos fallos debido a su gran dificultad para ser detectados por técnicas de control estadístico de procesos. En esta ocasión se logra detectar el fallo 14 que antes no se detectaba y que ahora lo conseguimos detectar en poco tiempo y además detectando la mayoría de las anomalías que este fallo produce en la planta. El fallo 19 no se llega a detectar porque no se supera el umbral 10 veces consecutivas pero se consiguen hallar más de un tercio de las anomalías producidas.

Por otra parte, si observamos la media de porcentaje de alarmas detectadas vemos que se ha logrado subir de un 53% a un 57% acercándose al porcentaje de alarmas detectadas por el anterior método utilizado (PCA). Sin



embargo, el tiempo de detección medio ha subido de 79 a 107. Es decir, se detectan más anomalías del sistema pero cuesta un poco más de tiempo confirmar si ha habido un fallo en la planta. Si realizamos la media sin los fallos 3,9 y 15 el porcentaje de alarmas detectadas sube hasta 66.20, un buen resultado que vamos a intentar mejorar con el siguiente paso.

Para mejorar estos datos, vamos a realizar otro cambio en la técnica que hemos utilizado. Seguiremos utilizando la distancia mahalanobis ya que nos está proporcionando buenos resultados. También vamos a seguir trabajando con el número de variables reducido a 33 como en el caso anterior ya que observamos que el porcentaje de alarmas detectadas subía considerablemente. El cambio que vamos a hacer en esta ocasión es cambiar las dimensiones a las que se reduce por defecto que estábamos trabajando con dos dimensiones por lo que se va a probar si reduciendo a 3 dimensiones se consiguen mejores resultados aunque se trabaje con un mayor número de datos.

Los resultados que se obtienen siguiendo la metodología anterior de cálculo son los siguientes:

Tipo de fallo	Tiempo de fallo	Alarmas detectadas (%)
IDV (1)	4	99,625
IDV (2)	13	98,75
IDV (3)	No detecta	4,5
IDV (4)	No detecta	4,5
IDV (5)	1	100
IDV (6)	1	100
IDV (7)	1	38,25
IDV (8)	21	96,125
IDV (9)	No detecta	4,375



IDV (10)	25	65,75
IDV (11)	302	17,875
IDV (12)	13	95,25
IDV (13)	45	94,75
IDV (14)	2	92,875
IDV (15)	685	10,5
IDV (16)	31	72,625
IDV (17)	25	84
IDV (18)	86	89,5
IDV (19)	No detecta	39,5
IDV (20)	80	54,25
IDV (21)	276	60,125
Media	<b>94,76</b>	<b>63,01</b>
Media sin los fallos 3,9 y 15	<b>57,88</b>	<b>72,43</b>

Tabla 11: Resultados obtenidos aplicando t-SNE utilizando 33 variables, la distancia mahalanobis y reduciendo a 3 dimensiones

En este caso, sí que se logra detectar uno de los fallos incipientes, en concreto el fallo 15, aunque tampoco se puede considerar un éxito ya que se detecta muy tarde y con muy poco porcentaje de alarmas detectadas. Se sigue sin detectar el fallo 19 pero se ha aumentado número de anomalías que se detectan llegando casi al 40%.

Por otro lado, si nos fijamos en los resultados generales que se han obtenido podemos observar que se ha conseguido mejorar de nuevo el número de alarmas detectadas llegando a un valor mucho más cercano al obtenido con el método PCA que con los anteriores métodos utilizando mediante la técnica t-SNE. También se ha reducido el tiempo que se tarda en detectar los fallos en comparación con el análisis anterior reduciendo los datos iniciales a un espacio de tres dimensiones en vez de a dos dimensiones.



En conclusión, después de haber realizado la técnica t-SNE probando diferentes parámetros como la distancia entre puntos utilizada o las dimensiones a las que se reduce el espacio original, hemos podido comprobar que los mejores resultados se han obtenido tras reducir el número de variables de la planta a 33, se han cambiado los datos que Matlab utiliza por defecto y se ha utilizado la distancia mahalanobis y una reducción dimensional a un espacio de  $\mathbb{R}^{96 \times 3}$ , consiguiendo un tiempo medio para detectar los fallos de 94,77 y un porcentaje de alarmas detectadas del 63,01. Si eliminamos los fallos 3,9 y 15 para obtener la media final, se observa que el tiempo de detección disminuye a 57,88 el mejor valor obtenido hasta ahora por el método t-SNE. Por otro lado, se obtiene un porcentaje de alarmas detectadas del 72,43% logrando también el valor más elevado con este método.

#### 4.4. Clasificación de las anomalías detectadas

Después de haber aplicado las dos técnicas de control estadísticos de procesos propuestas para lograr detectar las anomalías que producen los diferentes fallos de la planta química estudiada, se va a tratar de clasificar estas anomalías mediante la técnica de bosques aleatorios (*random forest*) y matrices de confusión. Las matrices de confusión son una herramienta para visualizar los datos de clasificación de fallos. Anteriormente estábamos trabajando con bloques de datos que correspondían en cada caso a un mismo tipo de fallo el cual sabíamos con anterioridad cual era en cada caso. Sin embargo, si se aplican las técnicas anteriores sin conocer el fallo previamente, podemos detectar que ha ocurrido una anomalía pero no saber cuál es la causa que lo ha producido. Por tanto, en este apartado se va a tratar de aplicar técnicas de clasificación utilizando ambos métodos de reducción de la dimensionalidad PCA y t-SNE, pudiendo comparar cual ofrece mejores resultados.



#### 4.4.1. Clasificación de fallos utilizando los datos de partida

En primer lugar, vamos a clasificar las anomalías detectadas utilizando las simulaciones de datos sin aplicar ningún método multivariante de reducción dimensional que hemos utilizado anteriormente.

Para aplicar la técnica de bosques aleatorios, primero es necesario tener dos matrices de datos  $X$  y  $X_t$ , siendo la primera la que se utiliza para crear el bosque y la segunda para poder testarlo y así experimentar si funciona correctamente. En nuestro caso, para la matriz  $X$  se ha utilizado el bloque de datos de 960 observaciones, sin embargo, en esta ocasión se han eliminado las 160 primeras ya que solo queremos los datos correspondiente a cuando hay un fallo presente en el sistema. Por otra parte, para la matriz  $X_t$  de test se ha utilizado el otro bloque de datos de 480 observaciones. En ambas matrices se han juntado todos los datos correspondientes a todos los fallos de la planta, sin embargo, se han eliminado los fallos 3, 9 y 15 que suelen dar problemas para ser detectados y por tanto, clasificados.

Se crea el bosque con la función `TreeBagger` de Matlab, en la que tenemos que fijar el número de árboles que queremos que tenga el bosque. Hemos probado con varios valores y se ha comprobado que el que mejor resultado da es utilizar un valor de 200 bosques de decisión. En esta función se utiliza la matriz  $X$  para crear el bosque. Después se utiliza el bosque con la función `predict` para obtener la matriz de salida  $T$  y  $T_{deseada}$  que nos permita finalmente crear la matriz de confusión para visualizar gráficamente los resultados. Las matrices de confusión obtenidas son las siguientes:





**Random forest original**

0	960	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	6.6%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
1	0	800	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
2	0	0	800	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
3	0	0	0	800	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
4	0	0	0	0	800	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
5	0	0	0	0	0	800	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
6	0	0	0	0	0	0	800	0	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
7	0	0	0	0	0	0	0	800	0	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
8	0	0	0	0	0	0	0	0	800	0	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
9	0	0	0	0	0	0	0	0	0	800	0	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
10	0	0	0	0	0	0	0	0	0	0	800	0	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
11	0	0	0	0	0	0	0	0	0	0	0	800	0	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
12	0	0	0	0	0	0	0	0	0	0	0	0	800	0	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
13	0	0	0	0	0	0	0	0	0	0	0	0	0	800	0	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	800	0	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	800	0	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	800	0	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	800	0	0	0	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%
	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%	100%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
		0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17			

Figura 26: Matriz de confusión utilizando los datos de creación del bosque

**Random forest test**

0	427	4	11	0	10	0	16	62	56	20	15	1	144	31	75	54	91	19	11.2%
	4.9%	0.0%	0.1%	0.0%	0.1%	0.0%	0.2%	0.7%	0.6%	0.2%	0.2%	0.0%	1.7%	0.4%	0.9%	0.6%	1.1%	0.2%	58.8%
1	0	473	0	0	0	0	29	0	0	0	0	0	0	0	0	0	0	0	94.2%
	0.0%	5.5%	0.0%	0.0%	0.0%	0.3%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	5.8%
2	0	0	453	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	99.8%
	0.0%	0.0%	5.2%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.2%
3	0	0	0	475	0	0	0	43	0	0	0	0	0	0	0	0	0	0	91.7%
	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	8.3%
4	3	0	0	0	402	0	4	0	10	12	0	0	0	2	1	2	0	0	92.2%
	0.0%	0.0%	0.0%	0.0%	4.6%	0.0%	0.0%	0.0%	0.1%	0.1%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	7.8%
5	0	0	0	0	0	478	0	0	2	0	0	0	0	0	0	0	0	0	99.6%
	0.0%	0.0%	0.0%	0.0%	0.0%	5.5%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.4%
6	0	0	15	0	12	0	307	0	0	16	48	0	0	0	0	0	0	0	7.1%
	0.0%	0.0%	0.2%	0.0%	0.1%	0.0%	3.5%	0.0%	0.0%	0.2%	0.6%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	12.9%
7	13	0	0	0	11	0	1	233	4	1	24	0	16	0	0	0	16	32	56.4%
	0.2%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	2.7%	0.0%	0.0%	0.3%	0.0%	0.2%	0.0%	0.0%	0.2%	0.4%	0.4%	13.6%
8	4	0	0	4	0	0	0	347	0	1	1	1	3	2	6	1	2	0	83.3%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	4.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.2%	0.4%	0.2%	5.7%
9	0	1	0	0	36	2	45	0	0	378	101	0	0	35	0	0	0	0	83.2%
	0.0%	0.0%	0.0%	0.0%	0.4%	0.0%	0.5%	0.0%	0.0%	4.4%	1.2%	0.0%	0.0%	0.4%	0.0%	0.0%	0.0%	0.0%	16.8%
10	0	0	0	0	0	0	62	0	0	19	269	0	0	16	0	0	0	0	73.5%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.7%	0.0%	0.0%	0.2%	3.1%	0.0%	0.0%	0.0%	0.2%	0.0%	0.0%	0.0%	0.0%	16.5%
11	0	0	0	1	0	0	0	0	18	0	0	0	459	0	0	0	0	0	95.6%
	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.2%	0.0%	0.0%	0.0%	5.3%	0.0%	0.0%	0.0%	0.0%	0.0%	4.4%
12	1	0	0	0	2	0	1	156	0	11	6	0	282	0	0	0	32	40	53.1%
	0.0%	0.0%	0.0%	0.0%	0.1%	0.0%	0.0%	1.8%	0.0%	0.1%	0.1%	0.0%	3.3%	0.0%	0.0%	0.4%	0.5%	0.6%	16.9%
13	0	0	0	0	1	0	0	5	3	0	17	2	432	0	4	7	0	0	91.7%
	0.0%	0.0%	0.0%	0.0%															



La arquitectura de la matrices de confusión como las de las figuras 26 y 27 siempre es la misma, las filas representan las clases que se predicen, es decir, los fallos del sistema predichos. Por otro lado, las columnas de la matriz muestran el conjunto de fallos deseados. Además, la traza diagonal se corresponde con los fallos correctamente clasificados, lo que significa que las observaciones que se encuentren fuera de la diagonal se corresponderán con fallos clasificados de manera errónea.

La última columna que se encuentra más a la derecha se corresponde con el número total de observaciones de cada fallo o clase que se ha predicho. Cuantificando con un porcentaje en verde la tasa de aciertos detectados (precisión) y en rojo la tasa de falsos descubrimientos. A su vez, la última fila muestra el total de observaciones para las clases deseadas, mostrando también en verde el porcentaje de verdaderos positivos (sensibilidad) y de falsos negativos. Finalmente, nos encontramos con la última celda de la diagonal en la esquina inferior derecha que nos aporta el porcentaje total de predicciones correctas reflejado por el color verde y de predicciones incorrectas en color rojo. Será el valor de más importancia para analizar que técnica clasifica mejor los fallos de la planta.

Una vez obtenido los resultados en las matrices de confusión, se pueden visualizar rápidamente cual son los resultados obtenidos. En la figura 26 probamos el bosque con los mismos datos utilizamos para su creación, por lo que es lógico que se identifiquen a la perfección todos los fallos. Observamos que no existe ningún fallo clasificado fuera de la diagonal, obteniendo una de predicciones correctas del 100%. Por otro lado, en la figura 27 nos encontramos con un resultado diferente debido a que estamos utilizando un bloque de datos distintos para testear el bosque. En este caso la tasa de predicciones correctas totales es del 79%, es decir, habría un 21% de anomalías que se han clasificado erróneamente.



#### 4.4.2. Clasificación de fallos utilizando PCA

En el primer apartado se ha utilizado la técnica de clasificación de fallos mediante el uso de los datos iniciales correspondientes a la medición de las 52 variables de la planta en cada instante. Por lo que en este apartado vamos a comprobar si podemos clasificar estos fallos reduciendo las 52 variables a sus componentes principales y así trabajar con un número reducido de datos.

Para ello utilizamos la misma metodología que en el primer caso. Se utilizarán dos matrices  $X$  y  $X_t$  una para crear el bosque y la otra para obtener los resultados de test. A diferencia que en el caso anterior, estas matrices tendrán una dimensión reducida ya que se aplicará el método de componentes principales PCA para reducir el número de variables. Para ello, simplemente calcularemos la matriz de correlación  $R$  (Ecuación 7) de la matriz inicial  $X$  y sus correspondientes valores singulares. Después, introduciremos la varianza con la que queremos trabajar que usaremos 99% un valor más elevado que el que se utilizó para detectar anomalías debido a que en este caso la matriz  $X$  contiene 14560 observaciones ya que estamos trabajando con todos los datos de fallos en una sola matriz. Finalmente se calculará la matriz  $P$  y obtendremos la matriz  $X$  final (Ecuación 9) con los componentes principales. Para la matriz  $X_t$  el procedimiento es mucho más sencillo ya que simplemente tendremos que multiplicar la matriz inicial de datos por la matriz  $P$  ya calculada para la matriz  $X$ .

A continuación se muestran los resultados obtenidos para un bosque de 200 árboles mediante las correspondientes matrices de confusión:

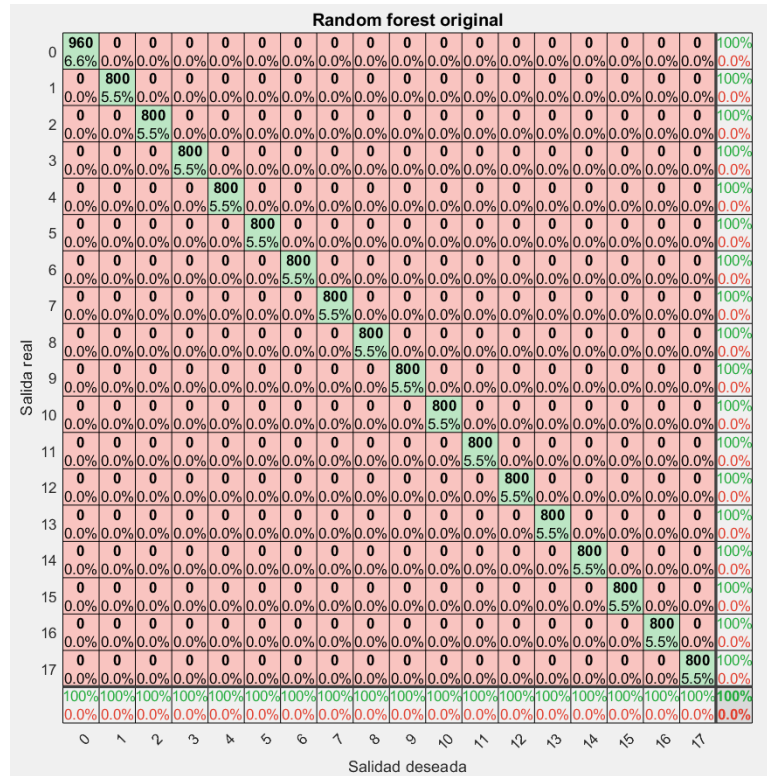
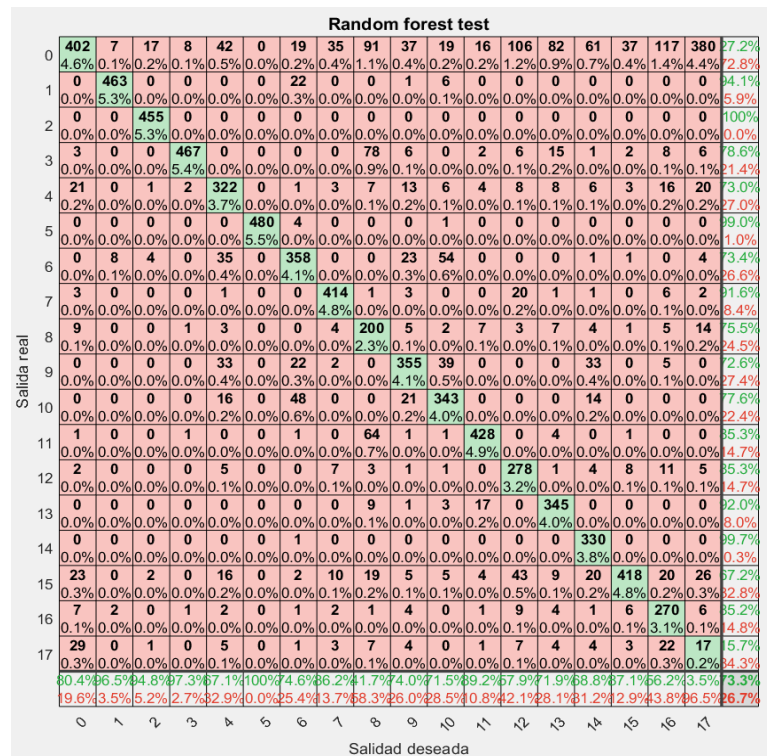


Figura 28: Matriz de confusión mediante PCA utilizando los datos de creación del bosque





En la figura 28 se observa que el bosque funciona perfectamente al utilizar los datos con los que se ha creado al igual que pasaba en el apartado anterior, obteniendo un 100% de fallos predichos correctamente. Sin embargo, lo interesante se encuentra en la figura 29 donde testeamos al bosque con otros datos distintos. Al igual que ocurría anteriormente, en este caso ya no se obtiene un porcentaje perfecto de anomalías clasificadas correctamente, obtenido un porcentaje total del 73,3% algo inferior al obtenido utilizando los datos sin reducción de dimensionalidad pero que sigue siendo bastante alto y fiable a la hora de clasificar los fallos producidos en la planta. Además, podemos observar como el último fallo estropea mucho la estadística ya que solo se clasifica correctamente 16 observaciones de las 480 que se disponen.

#### 4.4.3. Clasificación de fallos utilizando t-SNE

Finalmente, vamos a tratar de clasificar las anomalías detectadas mediante el uso de la técnica t-SNE. Al igual que en el apartado anterior, se reducirá la dimensión de los datos utilizados y observaremos si los resultados son similares a los ya obtenidos.

Para el cálculo de la matriz  $X$  que utilizaremos para crear el bosque, se van a implementar los parámetros que mejores resultados nos otorgaron a la hora de detectar las anomalías producidas en la planta. De esta manera, reduciremos las variables de 52 a 33, se utilizará la distancia mahalanobis y el número de dimensiones del nuevo espacio creado por t-SNE será tres. Al crear esta nueva matriz de datos  $X$  mediante la función *tsne* de Matlab, calcularemos la matriz  $A$  mediante la ecuación 26. Esta matriz es la que vamos a utilizar para calcular la matriz de test  $X_t$  a partir del otro bloque de datos distinto al utilizado para calcular la matriz  $X$ . Habiendo ya obtenido las matrices  $X$  y  $X_t$  necesarias, se procede a la creación del bosque con 200 árboles y al cálculo de las matrices de confusión, obteniendo los siguientes resultados:







En la figura 30 se observa que el método funciona de forma correcta si se utilizan los datos de creación del bosque, aportando el mismo resultado para esta prueba que los métodos anteriores. Sin embargo, obtener una tasa del 100% para estos datos no es relevante a la hora de analizar si el método clasifica correctamente las anomalías. Esto lo podemos comprobar mediante la figura 31, donde se ha utilizado la matriz  $X_t$  para testear el bosque y cómo podemos observar, los resultados son desastrosos ya que globalmente solo se clasifican correctamente 28,9% de las anomalías del sistema. Se observa que hay algunos datos como los correspondientes al fallo 1, 2 o 5 que se clasifican bastante bien pero el resto de los fallos son muy difíciles de clasificar a través de la técnica t-SNE.

Aunque la técnica t-SNE nos ayude reduciendo en gran medida los datos con los que se trabaja, no resulta eficiente utilizarla para diagnosticar y clasificar las anomalías que se detecten en el proceso. La razón de este resultado seguramente será el cálculo de la matriz de proyección A, para pasar del espacio de alta dimensión al de dimensión 3, que se ha considerado lineal, cuando la técnica t-SNE hace una reducción no lineal de la dimensionalidad. Para resolver este problema habrá que investigar en como calcular una nueva matriz de proyección que tendrá que ser no-lineal.



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos







# **CAPITULO V: CONCLUSIONES Y TRABAJO FUTURO**



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos





## 5.1. Comparación entre las técnicas de control estadístico de procesos utilizadas

La nueva industria 4.0 nos proporciona grandes cantidades de datos que aportan gran información sobre los procesos industriales que se llevan a cabo continuamente en todo el mundo. En este trabajo se ha buscado utilizar diferentes técnicas de control estadístico de procesos que nos permitan tratar todos esos datos obtenidos y lograr mejorar la calidad de un sistema. Para ello, se han tomado los datos de la planta Tennessee Eastman, un proceso estándar para trabajar con estas técnicas.

La primera técnica de análisis multivariante que se ha utilizado ha sido el método de Análisis de Componentes Principales (PCA), donde conseguimos reducir las 52 variables iniciales a 31 componentes principales con una varianza del 90% y un nivel de significancia de 0.02. Para ello se analizaron en primer lugar los datos de comportamiento sin fallo de la planta para crear los parámetros y umbrales necesarios para estudiar los datos con fallo del sistema. Se consiguieron detectar todos los fallos de la planta excepto los fallos 3, 9 y 15, muy difíciles de detectar por otros métodos. Por otro lado, si se consideran valores medios, la cantidad de datos detectados con las estadísticas  $T^2$  y  $Q$  son 62,90% y 77,64% respectivamente. Si nos fijamos en los valores obtenidos eliminando los datos correspondientes a los fallos 3, 9 y 15, se obtienen unos valores para las estadísticas  $T^2$  y  $Q$  de 69,48% y 86,75%, siendo este último el mejor resultado obtenido en cuanto a porcentaje de alarmas detectadas. También se consiguieron unos tiempos de detección medios para el estadístico  $T^2$  de 80,05 y para el estadístico  $Q$  de 97,6 que si eliminados los fallos 3, 9 y 15 se quedan en unos valores de 73,03 para  $T^2$  y 39,83 para  $Q$ . Se puede observar que al eliminar estos fallos, la media de tiempos baja, siendo el estadístico  $Q$  el que más rápido consigue detectar los fallos. También se pudo analizar cuál de las variables contribuía más a que se produjera un fallo en el sistema, donde observamos que las



variables que más solían fallar eran la 16 (Flujo de alimentación de A y C), 21 (Temperatura del reactor) y 46 (Concentración de F en la Purga).

Por otro lado, se utilizó la técnica no lineal de incrustación de vecinos estocásticos distribuidos en  $t$  (t-SNE). Se comenzó utilizando los valores predeterminados por Matlab para la función *tsne* y se aplicó la técnica en primer lugar a los datos sin fallo del sistema para conseguir los parámetros y umbrales de los estadísticos que son utilizados posteriormente para analizar los datos con fallo del proceso. Se obtuvieron buenos resultados que con el método PCA así que se probaron diversas modificaciones de los parámetros que utiliza el método como las distancias entre puntos o las dimensiones reducidas. Finalmente, los mejores resultados son proporcionados por la distancia mahalanobis con una reducción a tres dimensiones, además de una reducción de las variables del proceso. Los mejores datos obtenidos son un tiempo medio para detectar los fallos de 94,77 y un porcentaje de alarmas detectadas de 63,01%. Además, si nos fijamos en las medias sin los fallos 3,9 y 15 observamos que se obtienen un tiempo medio de detección de 57,88, un valor que supera el obtenido mediante el estadístico  $T^2$  con el método PCA. Por otra parte, el porcentaje de alarmas detectadas sube hasta 72,43%, superando también el obtenido por PCA con el estadístico  $T^2$ , pero sin llegar al valor obtenido mediante el estadístico Q.

Finalmente, se utilizó la técnica de árboles aleatorios para lograr diagnosticar las anomalías detectadas y poder predecir la causa que las origina. Se ha observado que el método de bosque aleatorio usando como datos de entrada los datos reducidos mediante la técnica t-SNE proporciona resultados pésimos para clasificar los datos analizados. Por otro lado, se obtuvieron buenos resultados para clasificar los fallos con el método de bosque aleatorio con los datos reducidos mediante PCA, aunque el mejor resultado se obtuvo al utilizar todos los datos del problema.



En conclusión, para detectar las anomalías del sistema se obtienen resultados similares para sendos métodos aunque los mejores datos se han conseguido mediante el uso de PCA y el estadístico Q. Sin embargo, el método PCA utiliza muchos más datos que t-SNE ya que la ventaja de este último es que se reduce a un espacio de dos o tres dimensiones, mientras que para PCA se utilizan 31 componentes principales aumentando en gran medida la cantidad de datos manipulados. Finalmente, mediante la técnica de árboles aleatorios usando los datos reducidos mediante la técnica t-SNE no se han podido clasificar de forma correcta las anomalías detectadas mientras que mediante los árboles aleatorios usando todos los datos o los datos reducidos mediante PCA se ha logrado fácilmente llegar a clasificar un 73,3% de las observaciones estudiadas.

## 5.2. Trabajo que se podría realizar en un futuro

Respecto a las mejores que se podrían hacer en un futuro, sería necesario buscar un método que permita clasificar los fallos mediante t-SNE para que lograra los mismos o mejores resultados que utilizando la técnica PCA ya que a pesar de que se logren detectar un porcentaje parecido de anomalías, es muy interesante poder saber su origen para poder arreglarlo y hacer que el sistema vuelva a funcionar correctamente sin fallo. También se podrían utilizar redes neuronales y diversos métodos de aprendizaje automático que nos permitan mejorar la clasificación de fallos y se complementen para ayudar también a la detección de anomalías del proceso. Por otra parte, se puede seguir estudiando diferentes parámetros de la técnica t-SNE que logren superar los resultados obtenidos mediante PCA. Además, se podrían obtener más datos de cada tipo de fallo para tener una visión más genérica de cómo funcionan ambos métodos ante grandes cantidades de datos.



Control de la calidad de un proceso mediante la detección y diagnóstico de anomalías  
usando técnicas de control estadístico de procesos





## Bibliografía

- [1] R. Costa, V. Puig y J. Blesa, «Introducción a la Diagnóstico de Fallos basada en Modelos mediante Aprendizaje basado en Proyectos,» *Revista Iberoamericana de Automática e Informática Industrial*, vol. 13, pp. 186-195, 2016.
- [2] K. Ishikawa, «¿Qué es el Control de Calidad?,» de *Introducción al control de calidad*, Díaz de Santos, 1989, pp. 13-59.
- [3] «Calidad Total: definición y modelos,» isotools, [En línea]. Available: <https://www.isotools.org/2015/05/01/calidad-total-definicion-y-modelos/>. [Último acceso: 17 Junio 2021].
- [4] C. Martínez, «¿Cuál es el origen y la utilidad de un sistema de gestión de calidad?,» *Revista digital Inesem*, 18 Agosto 2018.
- [5] Gladys, «Historia del SPC,» Measure Control, [En línea]. Available: <https://measurecontrol.com/la-historia-del-spc/>. [Último acceso: 20 Junio 2021].
- [6] M. J. de la Fuente Aparicio, «Apuntes del Control estadístico de procesos,» Universidad de Valladolid, 2019.
- [7] J. I. Contreras, «Altas Consultora,» 28 Enero 2021. [En línea]. Available: <https://www.atlasconsultora.com/como-controlar-la-variabilidad/>. [Último acceso: 20 Junio 2021].
- [8] I. Sánchez, «Introducción al control estadístico de procesos,» 2018. [En línea]. Available: [http://www.est.uc3m.es/esp/nueva\\_docencia/comp\\_col\\_leg/ing\\_tec\\_inf\\_gestion/estadistica/Documentacion/Temario\\_sinpres/ControlCalidad/Apuntes\\_Calidad.pdf](http://www.est.uc3m.es/esp/nueva_docencia/comp_col_leg/ing_tec_inf_gestion/estadistica/Documentacion/Temario_sinpres/ControlCalidad/Apuntes_Calidad.pdf). [Último acceso: 15 Julio 2021].
- [9] «SPC consulting group,» [En línea]. Available: <https://spcgroup.com.mx/>. [Último acceso: 21 Junio 2021].  
  
C. C. Florencio, «Gráficos de control basados en la variación de un rango multivariante,» 2016. [En línea]. Available: <http://hdl.handle.net/10016/26864>. [Último acceso: 15 Julio 2021].
- [11] L. v. d. Maaten, «Visualizing Data using t-SNE,» *Journal of Machine Learning Research*, vol. 9, pp. 2579-2605, 2008.
- [12] L. Kristjánssdóttir, «Exploración de grandes estudios espectroscópicos con métodos de aprendizaje automático,» 2018. [En línea]. Available: <https://riull.ull.es/xmlui/handle/915/10630>. [Último acceso: 15 Julio 2021].



- [13] «Sitio big data,» [En línea]. Available: <https://sitiobigdata.com/2019/10/27/una-introduccion-a-tsne-con-python/>. [Último acceso: 21 Junio 2021].
- [14] D. Liu, T. Gou y M. Chen, «Fault Detection Based on Modified t-SNE,» de *Symposium on Fault Detection, Supervision and Safety for Technical Processes (SAFEPROCESS)*, Xiamen, China, 2019.
- [15] V. Roman, «Aprendizaje Supervisado: Introducción a la Clasificación y Principales algoritmos,» Marzo 2019. [En línea]. Available: <https://medium.com/datos-y-ciencia/aprendizaje-supervisado-introducci%C3%B3n-a-la-clasificaci%C3%B3n-y-principales-algoritmos-dadee99c9407>. [Último acceso: 22 Junio 2021].
- [16] L. Breiman, «Random Forests,» *Machine Learning*, vol. 45, pp. 5-32, 2001.
- [17] V. Rodriguez, «Decision trees / árboles de decisión para clasificar en python,» 17 Octubre 2018. [En línea]. Available: <https://vincentblog.xyz/posts/decision-trees-arboles-de-decision-para-clasificar-en-python>. [Último acceso: 23 Junio 2021].
- [18] J. J. Downs y E. F. Vogel, «A plant-wide industrial process control problem,» *Compututer Chemical Engineering*, vol. 17, n° 3, pp. 245-255, 1993.
- [19] L. Chiang, E. Russell y R. Braatz, *Fault Detection and Diagnosis in Industrial Systems*, Springer-Verlag, 2000.