

**UNIVERSIDAD NACIONAL DE EDUCACIÓN A  
DISTANCIA**



**Departamento de Informática y Automática  
Escuela Técnica Superior de Ingeniería Informática**

**TÉCNICAS DE MINERÍA DE DATOS APLICADAS  
A FUSIÓN NUCLEAR: PREDICCIÓN EN TIEMPO  
REAL Y CLASIFICACIÓN**

**TESIS DOCTORAL**

**Giuseppe A. Rattá Gutiérrez**  
**Ingeniero Electrónico**

**2010**







UNIVERSIDAD NACIONAL DE EDUCACIÓN A  
DISTANCIA



**Departamento de Informática y Automática**  
**Escuela Técnica Superior de Ingeniería Informática**

**TÉCNICAS DE MINERÍA DE DATOS APLICADAS  
A FUSIÓN NUCLEAR: PREDICCIÓN EN TIEMPO  
REAL Y CLASIFICACIÓN**

**TESIS DOCTORAL**

**Autor: Giuseppe A. Rattá Gutiérrez, Ingeniero Electrónico**

**Director: Dr. Jesús Antonio Vega Sánchez**

**Tutor: Dr. Sebastián Dormido Canto**



**2010**









*A mi querida familia: mis padres Giuseppe y María del Pilar y mi hermana  
Carla, por el apoyo y el cariño que incondicionalmente me brindan.*



## Agradecimientos

La elaboración de una tesis doctoral es una tarea difícil. Se requieren cualidades limitadas: perseverancia, constancia, creatividad, dedicación. Sin embargo, de nada servirían estas cualidades si no existiese además un objetivo bien trazado y una guía adecuada para alcanzarlo. Es fundamental entonces darle gran parte (por no decir la mayor parte) del mérito a mi tutor, el **Dr. Jesús A. Vega**. Él ha sido un ejemplo, desde el primer momento, de trabajo, corrección, conocimientos y generosidad. A pesar de su preocupación por que siguiera el camino adecuado nunca dejó de fomentar las ideas que yo pudiera aportar, de discutir las conmigo y de motivarme. Su marca indudablemente permanecerá, y cada vez que en el futuro deba escribir un artículo o enfrentar un problema recordaré sus valiosos consejos.

Mis estancias en Inglaterra me dieron la oportunidad y el placer de trabajar con el **Dr. Andrea Murari**. Además de su reconocido nivel y conocimientos en casi todos los campos de la física de plasmas, su pasión por el trabajo y su curiosidad incesante fueron un estímulo gratificante. Sus constantes ideas resultaron ser un impulso decisivo para la solución de algunos de los problemas tratados.

Las colaboraciones con la UNED me permitieron conocer a mi tutor de tesis, el **Dr. Sebastián Dormido Canto**. Siempre lo recordaré con cariño como una persona de gran bondad e inteligencia.

El trabajo, para algunos, puede resultar aburrido e incluso tedioso. No ha sido mi caso. Parte de la alegría con la que me he levantado (descontando los 10 minutos posteriores al despertador, que no me producen ninguna alegría) cada día para venir al CIEMAT se la debo no sólo a la satisfacción que me produce la investigación sino también a mis compañeros. En el edificio 20 sería insuficiente llamarlos compañeros porque en realidad somos amigos. Atesoraré los momentos que pasé con cada uno de ellos, especialmente con los que por diversas circunstancias tuve el agrado de compartir más tiempo: **Guillermo Sánchez Burillo, Laura Barrera Orte, José Manuel García Regaña, Dr. Josep María Fontdecaba, Dr. Arturo Alonso, Dr. Iván Vargas, Dr. José A. Ferreira, Daniel Carralero y Dr. David Jiménez Rey**.

El apoyo de los **compañeros de la Unidad de Adquisición de Datos** fue constante desde que ingresé al centro. Me hicieron sentir como en casa y me ayudaron altruistamente cada vez que lo necesité.

Sería injusto olvidar a mis **familiares y amigos**. La estabilidad que ellos aportan, sobre todo cuando uno realiza estudios en un país extranjero, debe considerarse como un puntal en la vida y por lo tanto también en el trabajo.

El **grupo de física del Laboratorio Nacional de Fusión por Confinamiento Magnético** me ha demostrado que el paso del tiempo no disminuye el amor por la investigación sino que lo fortalece y profundiza.

A todos ellos y a los que en menor medida hicieron que esta experiencia fuese tan grata como edificante deseo agradecerles los momentos compartidos.

---

## COMENTARIOS INICIALES Y RESUMEN

El objetivo de este proyecto de Tesis es la aplicación de técnicas de minería de datos para la predicción en tiempo real y clasificación de eventos físicos en plasmas termonucleares.

La estructura del trabajo descrito en estas memorias incluye introducciones concisas sobre los principios de la fusión nuclear y de las técnicas de minería de datos aplicadas, detalladas en los capítulos 1 y 2 respectivamente.

En el capítulo 3, luego de revisar el estado del arte, se introducen métodos para la búsqueda y reconocimiento de patrones y de fenómenos físicos específicos en amplias bases de datos.

En el capítulo 4 se utilizan técnicas automáticas para interpretación de la información y la generación de modelos físicos a partir de los datos.

El capítulo 5 comienza con una rigurosa revisión del estado del arte sobre un fenómeno físico de gran importancia para la fusión nuclear: las disrupciones. Posteriormente se detalla el proceso de extracción de características realizado mediante la implementación de un considerable número de técnicas de análisis exploratorio. Estos métodos ayudan a determinar la correcta elección y preprocesado de las magnitudes físicas del plasma que evidencian comportamientos disruptivos.

En el capítulo 6 se profundiza el trabajo desarrollado en el capítulo 5, construyéndose un sistema de predicción de disrupciones capaz de trabajar en tiempo real y cuya original arquitectura le permite alcanzar las tasas de detección más altas hasta el momento. Este sistema está siendo instalado en el JET, el mayor Tokamak del mundo, situado en Culham (Reino Unido).

Finalmente, en el capítulo 7 se resumen y valoran los objetivos alcanzados y se proponen nuevas vías de trabajo futuro.



---

**ÍNDICE**

1	INTRODUCCIÓN A LA FUSIÓN NUCLEAR .....	11
1.1	INTRODUCCIÓN.....	11
1.2	LA FÍSICA DE PLASMAS.....	13
1.3	FUSIÓN POR CONFINAMIENTO MAGNÉTICO.....	16
1.3.1	CONFINAMIENTO EN DISPOSITIVOS TOROIDALES. ....	16
1.3.2	DIAGNÓSTICOS.....	19
1.4	EL “STELLARATOR” TJ-II. ....	19
1.4.1	INTRODUCCIÓN.....	19
1.4.2	CALENTAMIENTO.....	20
1.4.3	CONFINAMIENTO.....	22
1.5	EL TOKAMAK JET.....	23
1.5.1	INTRODUCCIÓN.....	23
1.5.2	CALENTAMIENTO.....	23
1.5.3	CONFINAMIENTO.....	25
1.5.4	DIAGNÓSTICOS.....	26
1.6	EL PROYECTO ITER. ....	27
1.6.1	DEMO. ....	28
2	TÉCNICAS DE MINERÍA DE DATOS APLICADAS A FUSIÓN NUCLEAR. 29	
2.1	INTRODUCCIÓN.....	29
2.2	EXTRACCIÓN DE CARACTERÍSTICAS.....	30
2.3	TÉCNICAS DE REDUCCIÓN DE DIMENSIONALIDAD Y VISUALIZACIÓN. ....	32
2.3.1	TÉCNICAS DE REDUCCIÓN DE DIMENSIONALIDAD.....	33
2.3.1.1	Análisis de Componentes Principales.....	33
2.3.1.2	Criterios para decidir el número de componentes. ....	34
2.3.1.3	Descomposición de valor singular (SVD). ....	34
2.3.2	ESCALADO MULTIDIMENSIONAL (MDS). ....	35
2.4	TÉCNICAS DE VISUALIZACIÓN. ....	37
2.4.1	TOURS VISUALES.....	37
2.4.2	MAPAS TOPOGRÁFICOS GENERATIVOS. ....	38
2.5	SISTEMAS DE APRENDIZAJE AUTOMÁTICO. ....	39
2.5.1	INTRODUCCIÓN.....	39
2.5.2	MÉTODOS DE APRENDIZAJE NO SUPERVISADO. ....	40
2.5.2.1	Introducción.....	40
2.5.2.2	K-medias o k-means. ....	41
2.5.1	MÉTODOS DE APRENDIZAJE SUPERVISADO. ....	43
2.5.1.1	Introducción.....	43
2.5.1.2	Árboles de clasificación y regresión.....	44
2.5.1.3	Máquinas de vectores soporte.....	45

---

3	RECONOCIMIENTO ESTRUCTURAL DE PATRONES EN GRANDES BASES DE DATOS .....	51
3.1	INTRODUCCIÓN Y ESTADO DEL ARTE.....	51
3.2	RECONOCIMIENTO DE PATRONES EN GRANDES BASES DE DATOS.....	53
3.2.1	BÚSQUEDAS EN BASES DE DATOS DE FUSIÓN: ESTADO DEL ARTE.....	53
3.2.2	BÚSQUEDAS DE FORMAS DE ONDA COMPLETAS.....	54
3.2.2.1	Extracción de características.....	56
3.2.2.2	Medida de similitud.....	56
3.2.3	INTRODUCCIÓN A LAS TÉCNICAS DE RECONOCIMIENTO MORFOLÓGICO PARA BÚSQUEDAS DE PORCIONES DE SEÑALES.....	57
3.2.3.1	Metodología general.....	57
3.2.3.2	Extracción de características.....	57
3.2.3.3	Medida de similitud.....	62
3.3	APLICACIONES DE TÉCNICAS ESTRUCTURALES EN JET.....	63
3.3.1	BÚSQUEDA GENERAL DE PATRONES.....	63
3.3.1.1	Extracción de características.....	63
3.3.1.2	Medida de similitud.....	65
3.3.2	BÚSQUEDA DE FENÓMENOS FÍSICOS ESPECÍFICOS.....	65
3.3.2.1	Introducción.....	65
3.3.2.2	Cortes en canales de ECE.....	65
3.3.2.3	Transiciones L-H.....	71
3.4	CONCLUSIONES.....	74
4	EXTRACCIÓN DE MODELOS FÍSICOS MEDIANTE TÉCNICAS DE APRENDIZAJE: LA TRANSICIÓN L-H.....	77
4.1.1	INTRODUCCIÓN.....	77
4.1.2	TRANSICIONES L-H.....	78
4.1.3	EXTRACCIÓN DE CARACTERÍSTICAS.....	79
4.1.4	BASE DE DATOS Y MODELOS TEÓRICOS.....	80
4.1.5	RESULTADOS.....	80
4.2	CONCLUSIONES.....	83
5	ANÁLISIS DIFERIDO DE DISRUPCIONES EN JET.....	85
5.1	INTRODUCCIÓN Y ESTADO DEL ARTE.....	85
5.1.1	FASES Y TIPOS DE DISRUPCIONES.....	88
5.1.1.1	Fases.....	88
5.1.1.2	Tipos de disrupciones.....	89
5.2	EXTRACCIÓN DE CARACTERÍSTICAS Y PROCESADO DE LAS DESCARGAS.....	91

---



---

5.2.1	BASE DE DATOS Y PRE-PROCESADO.....	91
5.2.2	SELECCIÓN DE LAS SEÑALES MÁS RELEVANTES Y USO DE VENTANAS TEMPORALES. ....	92
5.2.3	EXTRACCIÓN DE CARACTERÍSTICAS.....	95
5.2.3.1	Introducción.....	95
5.2.3.2	Vectores de características y sus representaciones visuales.....	95
5.3	IMPLEMENTACIÓN DE MÉTODOS SUPERVISADOS.....	99
5.3.1	VALIDACIÓN CRUZADA.....	99
5.3.2	RESULTADOS.....	100
5.4	CONCLUSIONES.....	102
6	PREDICCIÓN DE DISRUPCIONES EN TIEMPO REAL .....	103
6.1	INTRODUCCIÓN Y ESTADO DEL ARTE.....	103
6.2	BASE DE DATOS Y SIMULACIÓN DE TIEMPO REAL.....	108
6.2.1	BASE DE DATOS.....	108
6.2.1	SIMULACIÓN DE TIEMPO REAL.....	109
6.3	ARQUITECTURA DEL CLASIFICADOR Y ENTRENAMIENTO DE LOS MODELOS DE LA PRIMERA CAPA.....	111
6.4	FUNCIÓN DE DECISIÓN.....	116
6.4.1	ENTRENAMIENTO DE LA FD.....	117
6.4.1.1	La RID.....	117
6.4.1.2	Vectores de característica para el entrenamiento de la FD.....	118
6.5	PRUEBAS DE RENDIMIENTO.....	121
6.5.1	PRUEBAS INICIALES.....	122
6.5.2	PRUEBAS DEL SISTEMA PARA DIFERENTES CAMPAÑAS EXPERIMENTALES.....	122
6.6	COMPARACIÓN DEL RENDIMIENTO DEL SISTEMA PREDICTOR DESARROLLADO CON EL SISTEMA JPS.....	125
6.7	CONCLUSIONES.....	128
7	CONCLUSIONES.....	131
7.1	CONCLUSIONES.....	131
7.2	TRABAJOS FUTUROS.....	132
8	BIBLIOGRAFÍA.....	135
	APÉNDICE .....	145

---



---

## ÍNDICE DE TABLAS

Tabla 1.1 Parámetros generales del Tokamak JET.....	24
Tabla 3.1. Representación de dos descargas en corte encontradas, una de ellas en modo O y la otra en modo X. ....	69
Tabla 4.1 Los resultados de los diferentes modelos teóricos son comparados con los de las ecuaciones extraídas automáticamente del clasificador SVM desarrollado. ....	81
Tabla 5.1 Listado de las señales seleccionadas. ....	93
Tabla 5.2 Importancia relativa de cada señal (Relev.) asignada por el CART. ....	94
Tabla 6.1 Detalle de la base de datos recopilada para el análisis.....	109
Se dispara inmediatamente una alarma y los datos de interés son guardados para un posterior análisis en un formato semejante al representado en la Tabla 6.2. Se analiza la siguiente descarga volviendo al paso 1. ....	116
Tabla 6.3 Tres descargas en las que la función de decisión activó una alarma. ....	117
Tabla 6.4. Ejemplos de todos los tipos de resultados que pueden ser almacenados. ....	120
Tabla 6.5 Resultados de las pruebas para todas las secuencias. ....	123



## INDICE DE FIGURAS

Fig. 1.1 Energía de enlace por nucleón vs número másico para Fusión y Fisión nuclear. ....	13
Fig. 1.2 Representación de la reacción de fusión D-T.....	14
Fig. 1.3 Secciones eficaces vs temperatura para reacciones D-D, D-T y D- <sup>3</sup> He. ....	15
Fig. 1.4 Vista esquemática de un “Stellerator” (el TJ-II). ....	18
Fig. 1.5 Vista esquemática de un Tokamak. ....	19
Fig. 1.6 Activación del girotrón 2 en la descarga 15000 de TJ-II. ....	22
Fig. 1.7 Ejemplo esquemático de una configuración con limitador y una con punto X. ....	26
Fig. 2.1 Ejemplo de gráfico de autovalores. ....	34
Fig. 2.2. Ejemplo, en unidades arbitrarias, de una representación mediante GTM. ....	39
Fig. 2.3. Evolución del agrupamiento mediante K-means.....	42
Fig. 2.4 Esquema de aprendizaje de un sistema supervisado. ....	44
Fig. 2.5 Representación esquemática de un árbol de clasificación creado por el CART.....	46
Fig. 2.6 Ejemplo de separación lineal con SVM. ....	48
Fig. 2.7 Utilización de función Kernel. ....	50
Fig. 3.1 Ejemplo de esquema de clasificación jerárquico.....	55
Fig. 3.2 Asignación de letras con respecto a las pendientes de los segmentos de ajuste. ....	59
Fig. 3.3 Codificación “SID” de una señal. Todas las primitivas tienen una duración temporal idéntica. ..	60
Fig. 3.4 Codificación “SDV” de una señal ejemplo. ....	60
Fig. 3.5 Resultado de una búsqueda con primitivas SVD.....	62
Fig. 3.6 Esquema del protocolo de comunicación empleado en JET. ....	64
Fig. 3.7 Interfaz gráfica diseñada para la búsqueda de patrones dentro de señales. ....	64
Fig. 3.8 Señal de temperatura, correspondiente al canal 13, que entra y sale de un corte. ....	66
Fig. 3.9 Dos ejemplos de los resultados de las búsquedas de cortes en canales de temperatura. ....	69
Fig. 3.10 Identificación combinada para el reconocimiento de transiciones. ....	74
Fig. 5.1 Esquema del proceso de extracción de características. ....	96
Fig. 5.2 Visualizaciones bidimensionales de dos conjuntos de vectores de características (mediante PsGT y MTG). ....	97
Fig. 5.3 Representación de seis vectores de características. ....	100
Fig. 5.4 Porcentajes de aciertos de clasificadores SVM entrenados con diferentes características de ventanas disruptivas y no disruptivas.....	101
Fig. 6.1 La evolución temporal de cada descarga puede simplificarse como la concatenación de los vectores de características de 30 ms. ....	111
Fig. 6.2 Esquema del análisis en tiempo real mediante una estructura de clasificadores que actúan paralelamente en ventanas de tiempo consecutivas. ....	113
Fig. 6.3 Ejemplo gráfico de las ventanas seleccionadas en las descargas disruptivas y no disruptivas para el entrenamiento de cada uno de los modelos que componen la primera capa del predictor. ....	115
Fig. 6.4 Esquema general de entrenamiento de la función de decisión. ....	119

Fig. 6.5 Resultados de las pruebas para tres secuencias ..... 121

Fig. 6.6 Tasas de acierto del modelo previamente entrenado. .... 125

Fig. 6.7 Porcentajes acumulados de disrupciones detectadas para diferentes tiempos de alarma. .... 127

---

# **1 Introducción a la fusión nuclear**

---

## **1.1 INTRODUCCIÓN.**

Los requerimientos energéticos del hombre para cubrir necesidades de producción industrial, transporte, iluminación, climatización, investigación y servicios, se van incrementando año tras año. Estimaciones de la ONU indican que la población mundial alcanzará los 9400 millones de habitantes a mediados del siglo XXI. Este aumento inevitablemente generará un incremento de la demanda de energía. En los países desarrollados ya se están imponiendo criterios de ahorro energético y por lo tanto es previsible que la demanda de energía eléctrica no crezca al ritmo de la década anterior. Sin embargo, el consumo per cápita de los países en vías de desarrollo es diez veces inferior que el de los países industrialmente avanzados. A medida que estos países mejoren su nivel de vida su demanda energética también aumentará con rapidez. Según qué escenarios se considere para la evolución de la demanda energética, el consumo podría triplicarse a mediados del siglo XXI.

Los combustibles fósiles (carbón, petróleo, gas) son actualmente las principales fuentes de energía por su reducido coste. Estos recursos han sido masivamente explotados y las reservas existentes, principalmente de petróleo y gas, no garantizan su disponibilidad a largo plazo. Por otro lado, los residuos lanzados a la atmósfera durante su combustión tienen un impacto negativo en el medio ambiente, en especial los gases

de efecto invernadero y en particular el  $\text{CO}_2$ . Sin especular sobre si el calentamiento global es un hecho demostrado o simplemente un mito alarmista, en general se está de acuerdo en que existe una necesidad inminente de encontrar nuevas fuentes de energía que sean limpias, seguras y sostenibles.

En cuanto a las soluciones viables se destacan dos: las energías renovables y la energía nuclear. Las energías renovables son en general bien aceptadas por la sociedad al no plantear aparentemente riesgos para los trabajadores y ser consideradas amigables con el medioambiente. A pesar de ello, la base para la construcción de paneles solares, por ejemplo, para la producción de energía fotovoltaica es el silicio (basado en el procesado del sílice), cuya extracción es altamente costosa. Los molinos para la producción de energía eólica son poco rentables y su instalación, poco agradable a la vista. A nivel práctico, las energías renovables, en general, solamente pueden cubrir actualmente una pequeña fracción de los altos requerimientos energéticos demandados. El enfoque nuclear, por otro lado, se presenta como una solución poderosa ya que es capaz de producir grandes cantidades de energía. Esta puede ser generada mediante procesos de fisión nuclear, en las que se dividen núcleos de alto número másico, o de fusión nuclear, donde se juntan núcleos atómicos ligeros. La energía producida mediante fisión, descubierta en 1939 por los investigadores alemanes Hahn y Strassmann, se genera como una reacción en cadena. Este tipo de reacciones, hoy en día, puede ser controlada con seguridad. Sin embargo, la extracción, enriquecimiento y utilización de materiales de alto peso atómico (como el uranio) para los procesos de fisión nuclear generan residuos radioactivos que tardan miles de años en perder su nocividad para el medioambiente. Debido a estas circunstancias, existen ciertos reparos en parte de la opinión pública para la producción de energía mediante la fisión. La fusión termonuclear<sup>1</sup> en cambio, y en particular la obtenida mediante confinamiento magnético, requiere un preciso y continuo control del combustible para la operación del dispositivo y la pérdida de este control no conlleva peligros medioambientales. La energía obtenida mediante la unión de dos átomos ligeros es del orden de 4 veces mayor a la generada mediante la fisión, como puede observarse en la Fig. 1.1. Allí se representa la energía de enlace por nucleón frente al número másico  $A$ . Esta se obtiene

---

<sup>1</sup> La fusión no requiere de un tamaño crítico en el combustible sino que necesita una temperatura crítica. De allí el nombre de termonuclear.



dividiendo la energía de enlace del núcleo por sus  $A$  nucleones. Puede observarse en los núcleos livianos un aumento abrupto de la energía de enlace por nucleón frente al número másico.

Una ventaja añadida de la fusión es que no necesita una ardua tarea de extracción minera, ya que el Deuterio y el Tritio son abundantes en la naturaleza. Además, su obtención es significativamente más sencilla que la de elementos de alto número másico.

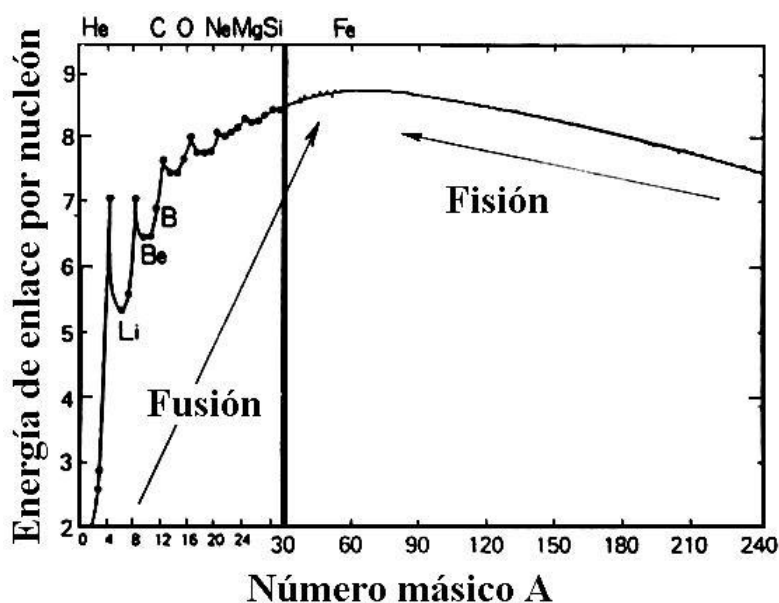


Fig. 1.1 Energía de enlace por nucleón vs número másico para Fusión y Fisión nuclear.

## 1.2 LA FÍSICA DE PLASMAS.

Desde los principios de la humanidad el hombre admiró la “inagotable” producción de energía generada por el sol, venerándolo en algunas civilizaciones como a un Dios. El primer paso científico hacia el desarrollo de una teoría que explicara esta poderosa fuente de radiación estelar se dio en 1905 cuando Einstein introdujo su famosa ecuación sobre la equivalencia entre masa y energía:  $E = mc^2$ , estableciendo que pequeñas variaciones en la masa pueden transformarse en una gran cantidad de energía. Francis Aston, 15 años después, logró otro avance fundamental al poder demostrarlo con resultados experimentales mediante la invención del espectrómetro de masas<sup>2</sup>.

<sup>2</sup> Instrumento que permite analizar con gran precisión la composición de diferentes elementos químicos e isótopos atómicos, separando los núcleos atómicos en función de su relación masa-carga.

Finalmente, en 1939 Hans Bethe publicó un artículo revolucionario explicando la generación de energía por fusión en las estrellas. El principio de la fusión es simple ya que involucra dos de los más abundantes elementos en el universo, el hidrógeno y el helio. La gran fuerza gravitacional que aparece en el interior de las estrellas (como en nuestro sol) comprime enormemente el hidrógeno presente y logra que se produzca la fusión de los núcleos. La energía generada es el resultado de la conversión de dos núcleos atómicos que se unen para formar uno de mayor masa atómica. Aplicando la fórmula de Einstein, esta diferencia de masa es transformada en energía (ver Fig. 1.2).

Sin embargo, las presiones y las temperaturas en el núcleo de las estrellas distan de las condiciones normales en la Tierra. En la década de 1950 se formalizaron los primeros intentos para la producción de energía mediante fusión [1]. Dentro de las posibles reacciones de fusión, las más atractivas son las denominadas como ciclo D-T (deuterio-tritio) por tener la mayor sección eficaz<sup>3</sup> a menores temperaturas [2], como puede observarse en la Fig. 1.3:

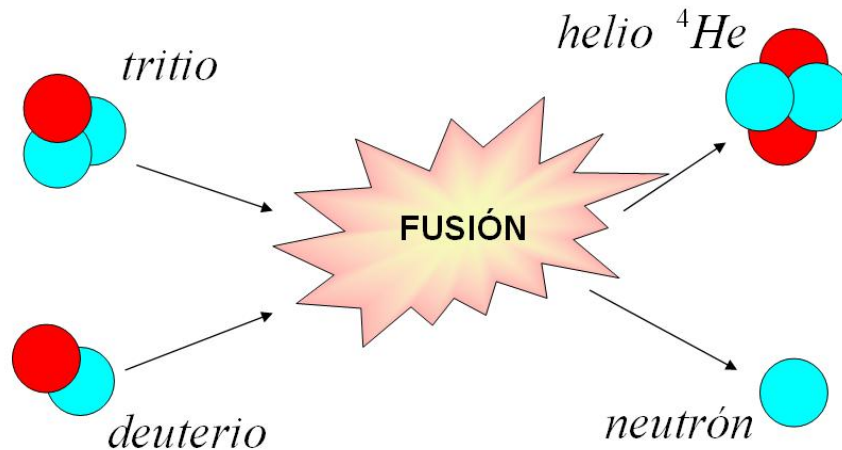
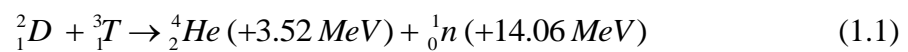


Fig. 1.2 Representación de la reacción de fusión D-T.

Un núcleo de deuterio se une a uno de tritio generando un núcleo de helio y un neutrón, liberando como consecuencia una energía de 17.58 MeV.

<sup>3</sup> Probabilidad de interacción entre dos partículas.

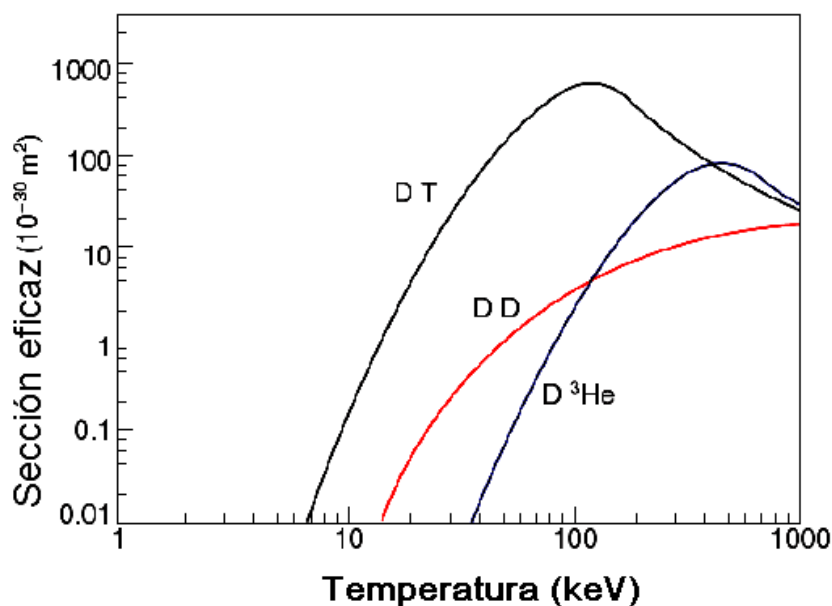
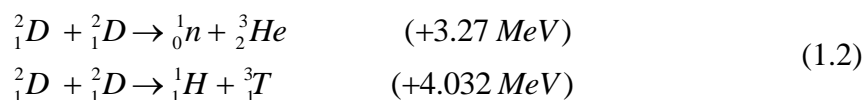


Fig. 1.3 Secciones eficaces vs temperatura para reacciones D-D, D-T y D-<sup>3</sup>He.

Como consecuencia de la fusión de los núcleos de deuterio y tritio se libera un total de 17.58 MeV. De ellos, 14.06 MeV se los lleva un neutrón en forma de energía cinética mientras los otros 3.52 MeV corresponden al núcleo de helio (o partícula  $\alpha$ ) generado. En plasmas D-T también pueden producirse reacciones D-D (deuterio-deuterio) de menores energías:



Como puede observarse, una de estas ramas genera tritio. El Deuterio es un elemento estable y abundante (30 gramos de cada metro cúbico de agua es deuterio). El tritio, sin embargo, no se encuentra en la naturaleza y debe producirse artificialmente. Puede generarse mediante la activación del hidrógeno contenido en el agua o por bombardeo de litio (elemento abundante en la corteza terrestre) con neutrones.

Los átomos de hidrógeno se encuentran totalmente ionizados a las temperaturas requeridas para que se produzcan las reacciones de fusión. A este estado de la materia se le denomina plasma. El plasma contiene un número significativo de partículas cargadas (iones y electrones) cuya dinámica presenta efectos colectivos [3].

Para que un reactor de fusión sea energéticamente rentable, las reacciones deberán generar una energía significativamente superior a la aplicada para la creación y el mantenimiento del plasma. Según el criterio de Lawson [4] la energía de las reacciones (en particular de las partículas  $\alpha$ ) resultan suficientes para calentar el plasma sin necesidad de un aporte energético externo. Lawson definió inicialmente un umbral mínimo, mediante el triple producto de la densidad electrónica ( $n_e$ ), la temperatura electrónica ( $\tau_E$ ) y el tiempo de confinamiento de la energía del plasma. Este triple producto, considerando un rendimiento del 33%, debe ser igual o superior a  $10^{21} keVs/m^3$  para que el plasma alcance la ignición<sup>4</sup>.

La investigación actual para el desarrollo de un futuro reactor de fusión se centra principalmente en dos vías, la fusión por confinamiento magnético y la fusión por confinamiento inercial. En los dispositivos de confinamiento magnético se utilizan campos para crear barreras que aislen el plasma de la cámara de vacío en la cual está contenido. El plasma, con una densidad de hidrógeno superior a  $10^{20} m^{-3}$ , se ha de calentar a temperaturas de decenas de KeV (siendo  $1eV \sim 11000 K$ ) con tiempos de confinamiento del orden de segundos. El confinamiento inercial, en cambio, se genera mediante el calentamiento y compresión de esferas de hidrógeno (las densidades alcanzadas son del orden de  $10^{23} m^{-3}$  y las temperaturas de decenas de KeV) durante períodos muy cortos (nanosegundos). Para la compresión y el calentamiento de las esferas se requieren láseres de alta potencia, cuyo desarrollo y posible aplicación en estado continuo plantean problemas de difícil solución [5].

### 1.3 FUSIÓN POR CONFINAMIENTO MAGNÉTICO.

#### 1.3.1 CONFINAMIENTO EN DISPOSITIVOS TOROIDALES.

Las primeras aproximaciones para conseguir un dispositivo de confinamiento magnético se basaron en solenoides que confinaban partículas radialmente permitiéndoles fluir axialmente. Para reducir las pérdidas, la solución propuesta se consistía en intensificar el campo magnético en los extremos (concepto conocido como *espejo magnético*). A pesar de dicho esfuerzo, el confinamiento resultante en este tipo

---

<sup>4</sup> Estado en el cual el calentamiento del plasma debido a las reacciones de fusión es suficiente para mantener su temperatura sin necesidad de fuentes externas de energía.

de dispositivos ha demostrado ser demasiado pobre por lo que han desaparecido del panorama de fusión.

Una solución alternativa para la reducción de pérdidas de partículas en los extremos es cerrar las líneas de campo magnético sobre sí mismas mediante la creación de campos magnéticos toroidales. En un sistema de campo magnético puramente toroidal, la curvatura y el gradiente del módulo de campo magnético provoca una deriva vertical en dirección opuesta para electrones e iones. Ésta crea un campo eléctrico que induce a su vez una deriva hacia el exterior en todo el plasma, por lo que una configuración puramente toroidal es intrínsecamente inestable. Para evitar esta separación de cargas es necesario retorcer las líneas de campo mediante una componente magnética adicional (poloidal). Una sola línea de campo describe así una superficie de flujo. El transporte<sup>5</sup> perpendicular al campo magnético  $B$  está restringido por la fuerza de Lorentz<sup>6</sup>, y por tanto, en una misma superficie de flujo (dirección paralela a las líneas de campo) los parámetros del plasma se pueden considerar prácticamente constantes, mientras que las principales variaciones se producen en la dirección perpendicular a las líneas de campo magnético.

Se destacan dos configuraciones en las que se pueden obtener reacciones de fusión por confinamiento magnético mediante estructuras de campo magnético con componentes toroidales y poloidales: las “Stellerator” y las Tokamaks.

El concepto de “Stellerator” fue ideado por el astrofísico Spitzer en 1951 [6]. La estructura de campo magnético en esta configuración se genera mediante conductores externos bobinados al toro (Fig. 1.4). La corriente que circula por las bobinas puede ser controlada desde el exterior y permite un modo de operación continuo. La intrincada geometría de las bobinas retuerce las líneas de campo magnético alrededor del toro. Sin embargo, la precisión de necesaria para el ensamblaje del dispositivo (en especial de las bobinas) es un verdadero reto para la construcción de este tipo de máquinas.

---

<sup>5</sup> Aquel proceso mediante el cual se produce una pérdida de partículas y energía desde el centro del plasma hacia el borde.

<sup>6</sup> Fuerza ejercida por el campo electromagnético que recibe una partícula cargada o una corriente eléctrica.

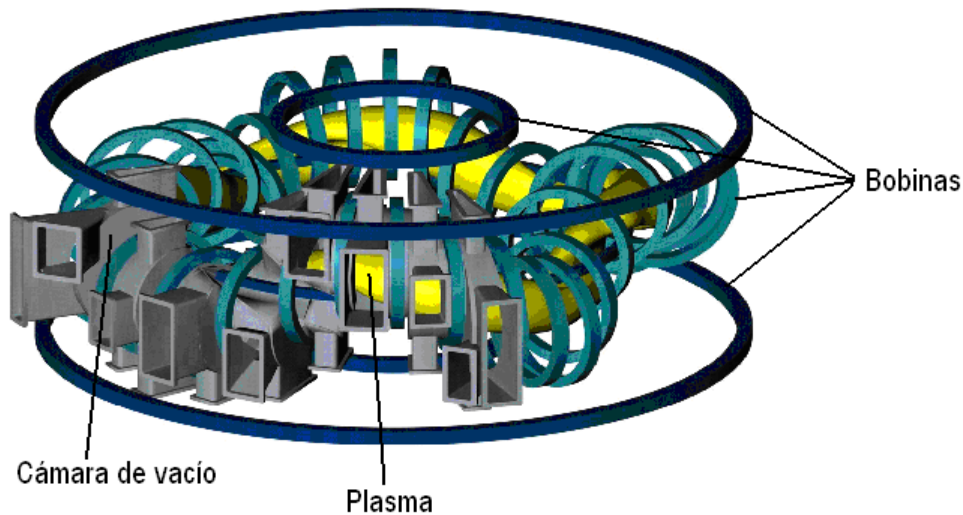


Fig. 1.4 Vista esquemática de un “Stellerator” (el TJ-II).

*Se representan el plasma (en amarillo), las bobinas (azul y celeste) y la cuarta parte de la cámara de vacío (gris).*

La configuración Tokamak (Fig. 1.5) fue propuesta por dos científicos rusos, Tamm y Sakharov. El nombre se deriva de las palabras rusas cuyo significado puede traducirse como “cámara toroidal con campo magnético”. La estructura de campo magnético se genera mediante bobinas toroidales y la componente poloidal la produce una corriente que se hace circular por el plasma. Esta corriente se induce mediante un transformador central. Hoy en día la configuración Tokamak es la más difundida y la que se utilizará para el futuro dispositivo experimental ITER (“International Thermonuclear Experimental Reactor”). Desde el punto de vista de ingeniería, presenta menos complicaciones para su construcción que un “Stellerator”. Sin embargo, una de sus principales desventajas se debe a pérdidas abruptas e inevitables de la energía del plasma (que ocurren sólo en este tipo de configuraciones), llamadas disrupciones.

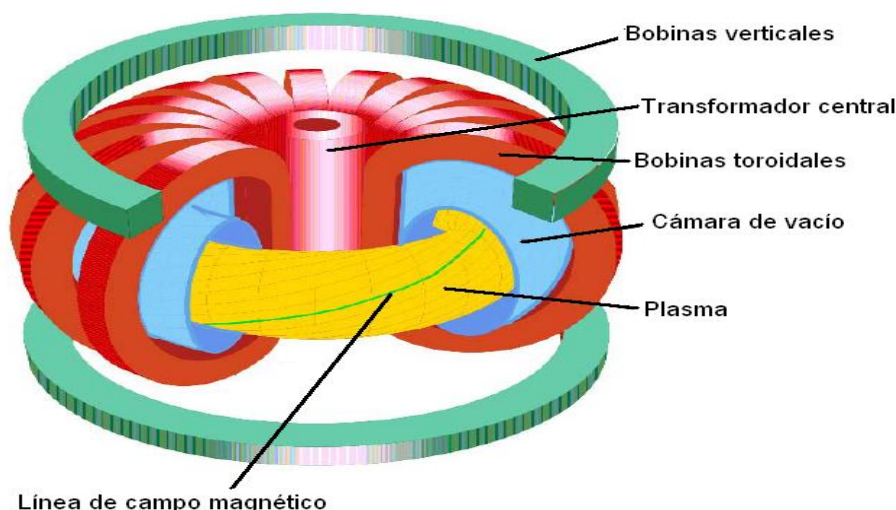


Fig. 1.5 Vista esquemática de un Tokamak.

*La corriente que circula por el plasma es inducida mediante el transformador central. Esta corriente aporta la componente poloidal del campo magnético.*

### 1.3.2 DIAGNÓSTICOS.

Los dispositivos desarrollados hasta el momento tienen como principales objetivos el estudio de los fenómenos físicos que ocurren en el plasma así como acercarse a las condiciones de ignición. El seguimiento de la evolución del plasma no es una tarea sencilla. Tal como se ha mencionado con anterioridad, el plasma puede encontrarse a temperaturas de decenas de KeV (unas diez mil veces la temperatura de la superficie del sol) y está confinado dentro de una cámara de vacío a cuyo interior es difícil acceder. Para la inspección del plasma y su evolución se requieren sistemas especiales de medida, llamados diagnósticos en la nomenclatura habitual de fusión. Los sistemas de diagnóstico son instrumentos especialmente desarrollados o adaptados para este tipo de complejas mediciones. Cada dispositivo para el estudio de la fusión dispone de un conjunto de diagnósticos que posibiliten determinar las características principales de los plasmas producidos.

## 1.4 EL “STELLARATOR” TJ-II.

### 1.4.1 INTRODUCCIÓN.

En Madrid se encuentra en funcionamiento el “Stellarator” más grande de Europa, el TJ-II [7]. Es de tipo heliac flexible [8] de tamaño medio. Sus objetivos principales son el estudio del confinamiento y del transporte (de energía y de partículas)

en plasmas calientes. El diseño de sus bobinas permite una gran variedad de configuraciones magnéticas, y un amplio rango de valores de la transformada rotacional<sup>7</sup>. Esta puede variar, en el centro del TJ-II, entre 0.9 y 2. Esta característica lo distingue de otros Stellarators al ser capaz de generar una gran variedad de geometrías magnéticas del plasma. Su eje magnético es helicoidal, confiriéndole una extremada tridimensionalidad. El TJ-II tiene una periodicidad  $M=4$  con un campo magnético central de  $\sim 1$  T y un radio mayor medio de 1.5 m.

Los estudios iniciales sobre la configuración magnética del TJ-II se realizaron conjuntamente entre el laboratorio ORNL (Oak Ridge, EE.UU.) y el CIEMAT [9]. La máquina, incluyendo la cámara de vacío, las bobinas y la estructura que las soporta, tiene un diámetro de cinco metros, una altura de dos metros sobre la base de la plataforma experimental y un peso de sesenta toneladas. Las primeras descargas con plasma se consiguieron a finales de 1997 [10]. Desde entonces, y como en todas las máquinas para el estudio de la fusión por confinamiento magnético, la operación del dispositivo se organiza en campañas experimentales que duran varios meses.

#### 1.4.2 CALENTAMIENTO.

El sistema más utilizado para calentamiento de plasmas en Stellarators está basado en la Emisión Ciclotrónica Electrónica Resonante, ECRH (“Electron Cyclotron Resonance Heating”). Esta emplea ondas electromagnéticas sintonizadas al segundo armónico del giro ciclotrónico de los electrones del plasma alrededor de las líneas de campo [11]. Las ondas son inyectadas perpendicularmente a las superficies magnéticas de la configuración. La frecuencia a la que son absorbidas las ondas depende de la intensidad del campo magnético, existiendo un valor máximo de la densidad electrónica para la cual no se propaga ECRH (fenómeno llamado de “reflexión total”). Cuando se alcanza este nivel de densidad las ondas son reflejadas y no calientan el plasma. Para los parámetros de las microondas inyectadas en TJ-II la densidad central de corte es de  $1.2 \cdot 10^{19} \text{ m}^{-3}$ .

El TJ-II tiene instalados dos girotrones (generadores de microondas de alta potencia), sintonizados al segundo armónico de la resonancia ciclotrónica electrónica.

---

<sup>7</sup> Rotación de la línea de campo, en radianes, por cada tránsito toroidal.



Cada girotrón es capaz de inyectar en el plasma 300 kW a la frecuencia de 53.2 GHz durante un pulso de duración máxima de 1 segundo.

El segundo sistema de calentamiento está compuesto de dos inyectores de haces neutros (NBI, del inglés “Neutral Beam Injection”). Este sistema consiste en la inyección de átomos neutros en el plasma pre-calentado. El haz de átomos posee una alta energía cinética que al colisionar con las partículas del plasma (entre ellas electrones) se cargan eléctricamente y consecuentemente quedan atrapadas por el campo magnético del dispositivo. Estos nuevos iones poseen una velocidad mucho mayor que la media de las partículas del plasma y consecuentemente provocan una serie de colisiones: ión-ión, ión-electrón y electrón-electrón. La velocidad media de los átomos confinados se incrementa en consecuencia. Mediante los dos inyectores (NBI-1 co-inyección y NBI-2 contra-inyección) se introducen en el plasma haces de átomos de hidrógeno de alta energía (31 keV) en dirección paralela al campo magnético toroidal. Los átomos de hidrógeno inyectados entran en la cámara de vacío con una energía suficiente como para llegar al borde del plasma sin desviarse para finalmente alcanzar su centro. Como resultado de las colisiones elásticas e inelásticas con los electrones e iones del plasma, los átomos del haz pueden ionizarse y ceder o no parte de su energía. La eficiencia de este tipo de calentamiento depende en gran medida de la densidad de partículas existente previamente, pero también de la distribución espacial de los campos magnéticos y eléctricos que encuentra el haz.

Para alimentar durante una descarga todo el sistema eléctrico de alta potencia del TJ-II (el cual incluye las bobinas, los sistemas de calentamiento y determinados diagnósticos) se dispone de un generador impulsional de 140 MVA, que almacena 100 MJ a 15 kV y 100 Hz. La duración de los pulsos eléctricos en el dispositivo TJ-II es de aproximadamente unos dos segundos, en los cuales se alimentan las corrientes de las bobinas, se introducen los gases dentro de la cámara de vacío y se encienden los girotrones para que calienten el gas. Los girotrones comienzan típicamente en 1020 ms, cuando termina de alcanzarse un valor estacionario en las corrientes de las bobinas (“*plateau*”) y se mantienen emitiendo microondas durante un intervalo típico de 250 ms, pudiendo alcanzar un máximo de 300 ms (véase la Fig. 1.6). Una vez que se apagan, el plasma se enfría y pierde su energía. El haz de neutros puede ser inyectado en diferentes instantes temporales, dependiendo del experimento que se vaya a realizar.

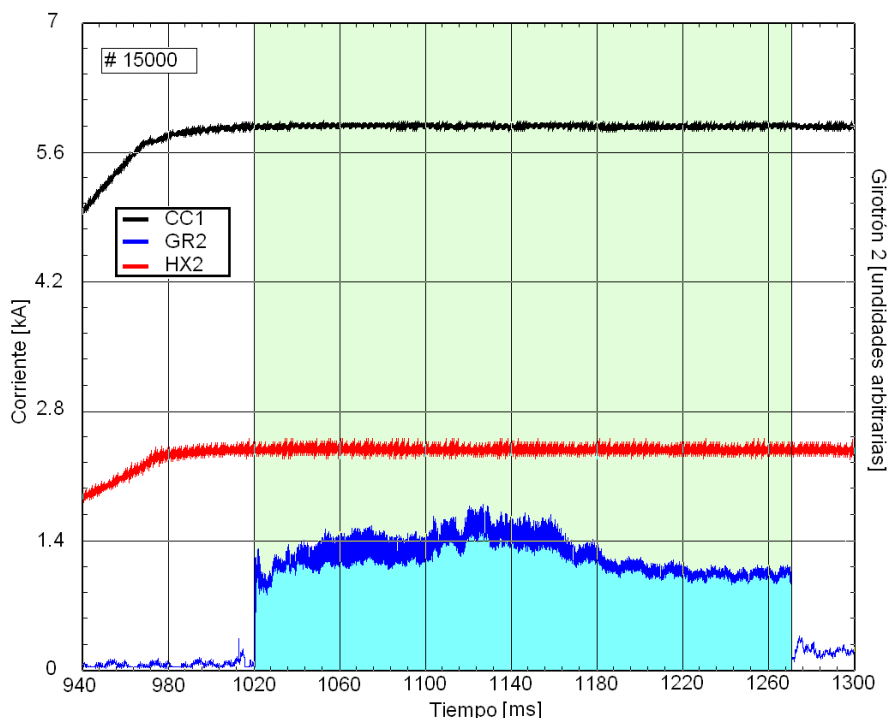


Fig. 1.6 Activación del girotrón 2 en la descarga 15000 de TJ-II.

Tanto para la corriente correspondiente a la bobina central (CC1) como a una de las helicoidales (HX2), el “plateau” termina de alcanzarse alrededor de los 1020 ms, momento en el cual suelen activarse los girotrones durante unos ~250 ms.

### 1.4.3 CONFINAMIENTO.

El campo magnético del TJ-II es generado por un conjunto de bobinas externas, que le confieren la configuración deseada. Cada bobina posee un número variable de vueltas, cuya corriente por vuelta depende del tipo de bobina y de la configuración magnética que se desee, llegando desde los 5.2 kA por vuelta en las bobinas radiales hasta los 32.5 kA por vuelta en las toroidales. La magnitud del campo magnético producido puede alcanzar ~1 T. La zona de corriente estabilizada tiene una duración nominal, para todas las bobinas, de 1 s con una frecuencia de repetición máxima de un pulso cada cinco minutos. Todas las bobinas se refrigeran con agua que circula a través de orificios longitudinales en los conductores de cobre y están embridadas por una estructura mecánica para evitar deformaciones.

La cámara de vacío se compone por ocho octantes iguales diseñados para aprovechar la simetría de los campos magnéticos. Cuenta con 88 ventanas de observación para los diagnósticos del plasma, los sistemas de calentamiento y los de inyección de gas. Con el fin de reducir el efecto de las impurezas pesadas en el plasma,

la cara interior de la cámara de vacío se recubre de boro o litio. Con estas técnicas se pueden conseguir descargas con baja concentración de impurezas y apropiado control de densidad [12].

Las líneas de campo magnético forman las superficies magnéticas, que son superficies cerradas y anidadas dentro de la cámara de vacío. El movimiento de las partículas cargadas está restringido por estas superficies [13].

## 1.5 EL TOKAMAK JET.

### 1.5.1 INTRODUCCIÓN.

JET (“*Joint European Torus*”) es el dispositivo de fusión por confinamiento magnético más grande del mundo. Su operación comenzó en 1983 y fue la primera instalación capaz de producir una cantidad significativa de energía mediante fusión controlada (cerca de 2 MW) en la campaña de deuterio-tritio de 1991 [14]. Consta de 32 bobinas en forma de D que se encuentran equi-espaciadas alrededor de la máquina. Estas bobinas son utilizadas para generar la componente toroidal de campo magnético. El transformador utilizado para inducir la corriente del plasma y a su vez generar la componente poloidal del campo está situado en el centro del dispositivo (Fig. 1.5).

Debido a la necesidad de un transformador para generar la corriente del plasma, la operación de JET (y la de todos los Tokamaks) es pulsada. Estos pulsos pueden generarse en JET con una frecuencia de uno cada veinte minutos y cada uno de ellos puede durar hasta sesenta segundos [15]. El plasma queda confinado en una cámara de vacío en forma de toroide cuyo radio mayor es de 2.96 metros, teniendo las secciones transversales unas dimensiones alto/ancho de 4.2 metros y 2.5 metros respectivamente. Los parámetros generales del diseño original de JET se detallan en la Tabla 1.1.

### 1.5.2 CALENTAMIENTO.

Además del calentamiento óhmico<sup>8</sup> y del NBI, el dispositivo consta de otros sistemas de calentamiento: Radio Frecuencia Ciclotrónica Iónica (ICRF, del inglés Ion

---

<sup>8</sup> El calentamiento óhmico del plasma es aquel producido por el paso de la corriente a través de él (efecto Joule). Al aumentar la temperatura el plasma pierde resistividad y como consecuencia el calentamiento óhmico pierde eficiencia.

Cyclotron Radio Frequency) y Corriente Híbrida Menor (LHCD, del inglés Lower Hybrid Current Drive).

Radio mayor del plasma	2.96m
Radio menor del plasma	2.10m(vert)-1.25m(horiz)
Longitud temporal de pulso (flat-top <sup>9</sup> )	20s
Peso del núcleo de hierro	2800t
Potencia de las bobinas de campo Toroidal	380MW
Campo magnético toroidal (en el eje del plasma)	3.4T
Corriente del plasma	3.2MA (Plasma circular) 4.8MA (Plasma con forma D)
Potencia adicional de calentamiento	25MW

*Tabla 1.1 Parámetros generales del Tokamak JET.*

Aunque el sistema LHCD puede tener un efecto de calentamiento ineficiente, es útil al generar en el plasma corrientes de varios MA. Las ondas electromagnéticas producen frecuencias resonantes mediante la modulación de un haz electrónico. En JET existen 24 tubos klystrons instalados en 6 módulos independientes. La onda electromagnética se transmite a la antena LHCD por un complejo sistema de guías de onda. Estas guías de onda consisten en conductores metálicos huecos con secciones del mismo tamaño que la amplitud de la onda transmitida. La antena debe ser montada directamente en la pared interior del dispositivo y estar tan cerca como sea posible del plasma.

El calentamiento por ICRF se utiliza en la mayoría de los experimentos de JET. Es resonante con la segunda frecuencia armónica de giro de los iones más comunes en plasmas de JET (deuterio) o con la frecuencia de giro de otras especies (como el tritio o

---

<sup>9</sup> Período a lo largo de la duración de la descarga que se caracteriza por corrientes y temperaturas máximas y estables.

el helio). La potencia total del sistema ICRF de JET alcanza los 32 MW aunque en la práctica suele aplicarse sólo una parte de ella. Estas potencias son enormes si se las compara, por ejemplo, con una emisora de televisión, donde un transmisor de 50 kW ya es considerado más que suficiente. Las líneas de transmisión que conducen las ondas desde los generadores hasta el Tokamak son cables coaxiales de bajas pérdidas. Estos cables son conductores en forma de tubo metálico que contienen en el centro, aislado, otro conductor. Este tipo de cables se utiliza en la mayoría de las transmisiones a altas frecuencias, como en el caso de las señales satelitales de televisión que son capturadas por una antena parabólica y que mediante cables coaxiales conducen las señales a un televisor. Sin embargo en JET, debido a las altas potencias del sistema, ha sido necesario ensanchar cientos de metros de estos cables hasta los 20 cm de diámetro, confiriéndoles una apariencia de tubería. Las líneas de transmisión terminan en 4 antenas ICRF instaladas en la pared interior de la máquina. Las ondas electromagnéticas ICRF no se pueden propagar a través de la cámara de vacío, ya que su longitud de onda es demasiado grande, por lo que la antena debe ser ubicada lo más cerca posible del plasma.

### 1.5.3 CONFINAMIENTO.

Como se explicó anteriormente, en todos los Tokamaks la componente poloidal del campo magnético se genera mediante la corriente del plasma y la toroidal mediante bobinas externas. La superficie magnética puede delimitarse o bien mediante una estructura sólida (placas llamadas limitadores) o estar completamente definidas mediante una configuración llamada punto X (*X-point configuration*), tal como puede apreciarse en la Fig. 1.7. El principal inconveniente de los limitadores es que son una fuente de impurezas que influyen negativamente en el plasma. Por otro lado, en plasmas elongados se genera una superficie magnética particular llamada separatriz (*separatrix*) que divide dos topologías magnéticas diferentes. Dentro de la separatriz se forman las superficies magnéticas cerradas y anidadas requeridas para un buen confinamiento. Fuera de la separatriz las superficies se encuentran abiertas. Una configuración con punto X es más estable comparada con las resultantes del uso de limitadores pero el volumen del plasma obtenido es menor. En general está aceptado que la configuración de punto X es la más adecuada para los Tokamaks y es la que actualmente se utiliza en JET.

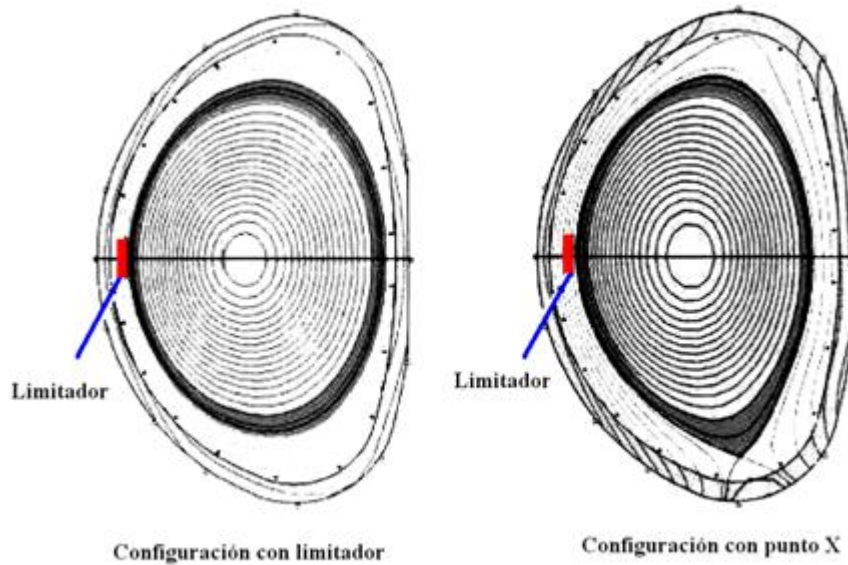


Fig. 1.7 Ejemplo esquemático de una configuración con limitador y una con punto X.

En la configuración con punto X la última superficie magnética no alcanza a entrar en contacto con el limitador.

#### 1.5.4 DIAGNÓSTICOS.

Actualmente en JET se utilizan unos 80 sistemas de diagnósticos. Además, existen alrededor de 20 más que en proceso de construcción. Algunos de los más comúnmente utilizados en JET se resumen en la siguiente lista:

- Bobinas magnéticas: proveen mediciones de campo magnético (incluyendo la detección de modos magnetohidrodinámicos), corriente, inductancia y energía del plasma.
- Dispersión Thompson: miden temperatura y densidad electrónica, así como perfiles de densidad.
- Cámaras ultra-rápidas de luz visible.
- Sistemas de interferometría: miden densidad integrada de línea.
- Antenas ECE (Electron Cyclotron Emission): miden temperatura y densidad electrónica.
- Espectrómetros de luz visible, ultravioleta y rayos X: miden temperaturas y densidades.

- Bolómetros: proveen mediciones de pérdidas de energía por radiación en el plasma.
- Sondos: son insertadas en el plasma para tomar mediciones directas.
- Cámaras de rayos X blandos: útiles para examinar propiedades magnéticas del plasma.
- Monitores de rayos X.

## 1.6 EL PROYECTO ITER.

El proyecto ITER (International Thermonuclear Experimental Reactor) es un consorcio internacional formado para demostrar la viabilidad científica y tecnológica de la fusión nuclear. Los participantes son la Unión Europea (UE), la Federación Rusa (en reemplazo de la Unión Soviética), Estados Unidos (aunque entre 1999-2003 decidió no participar), Japón, China (desde febrero 2003), Corea del Sur (desde mayo 2003) e India (desde diciembre 2005). Entre 1992-2004 participó Canadá. El principal objetivo de ITER es demostrar que es posible mantener el control de la combustión de un plasma de deuterio-tritio en condiciones próximas a ignición en estado estacionario. Deberá verificarse la validez de:

- Los componentes del reactor tales como los sistemas de extracción de energía y partículas del plasma.
- La disponibilidad e integración de las tecnologías esenciales en un reactor de fusión, como las bobinas superconductoras y el mantenimiento remoto.
- Los módulos de generación de tritio, la extracción de calor y generación de electricidad.

ITER será un dispositivo que integrará todos los aspectos físicos y tecnológicos necesarios para la producción de energía (menos los sistemas de producción de electricidad). La fase de diseño de ITER-EDA (Engineering Design Activities) preveía la producción de 1500 MW de potencia de fusión en un pulso de unos 1000 segundos. Para alcanzar la ignición se estimaban unos 100 MW de potencia de calentamiento adicional. La revisión ITER-FEAT (Fusion Energy Advanced Tokamak) tiene como finalidad la construcción de un dispositivo con un menor coste a cambio de unos

objetivos más modestos, previendo alcanzar una potencia de fusión de unos 500 MW con una razón entre la potencia de fusión y la potencia de calentamiento auxiliar ( $Q$ ) por lo menos de 10 en funcionamiento inductivo y de 5 en plasmas con corriente no inductiva, con una duración de la descarga suficientemente larga para demostrar la operación en estado estacionario (del orden de 1000 segundos).

#### 1.6.1 DEMO.

A largo plazo se proyecta realizar prototipos de reactores de fusión con el fin de la construcción de centrales comerciales que respondan a la demanda de la sociedad. De esta manera, después de ITER, se prevé la construcción del reactor de demostración DEMO, el cual producirá electricidad y deberá demostrar la viabilidad económica de una central de fusión. Se estima que para tal fin DEMO deberá superar en un 15% las dimensiones de ITER y conseguir densidades un 30% más altas.

Sin embargo existe todavía un largo camino antes hasta llegar a DEMO. En este deberá mejorarse tanto la tecnología e ingeniería de los dispositivos, como también alcanzar un conocimiento más profundo del comportamiento del plasma en diferentes condiciones de operación.



---

## 2 Técnicas de minería de datos aplicadas a fusión nuclear

### 2.1 INTRODUCCIÓN.

La minería de datos es un término que engloba una amplísima gama de técnicas de análisis de datos y de extracción y creación de modelos. Es posible aplicarlas a conjuntos de datos para describir tendencias, predecir comportamientos, extraer patrones o características, realizar búsquedas e interpretar la información. Es útil para el mejor aprovechamiento de los datos con el fin de extraer conocimientos de ellos. La minería de datos incluye técnicas tan tradicionales como las basadas en el análisis estadístico o tan modernas como los sistemas de aprendizaje. Las disciplinas de las que se nutre pueden resumirse en los siguientes grupos:

- **Las orientadas a las bases de datos y a la recuperación de la información:** técnicas de indexación y de búsquedas heurísticas a partir de conceptos o palabras claves, lo que puede considerarse como un proceso de clasificación de la información guardada para su mejor acceso.
- **Las estadísticas:** universalmente aceptadas para la interpretación de los datos dependiendo de su distribución, valores medios, varianza, técnicas bayesianas, modelización paramétrica y no paramétrica, etc.
- **Los sistemas de aprendizaje automático:** técnicas denominadas como “Inteligencia Artificial” que consisten en sistemas computacionales capaces de

aprender mediante ejemplos. Una vez que los sistemas han desarrollado modelos en el entrenamiento son capaces de resolver dicho problema.

- **Las técnicas para la visualización de datos:** útiles para la extracción de características y la comprensión e interpretación de información multidimensional.
- **Otras:** dependiendo de problemas específicos a resolver, cualquier combinación de las técnicas mencionadas se puede considerar como una técnica híbrida de minería de datos.

En los estudios realizados en esta Tesis se utilizaron métodos de todos los grupos mencionados anteriormente. Estos serán descritos brevemente en las próximas secciones.

## 2.2 EXTRACCIÓN DE CARACTERÍSTICAS.

La extracción de características consiste en identificar propiedades en cantidades medibles (que definirán un objeto) de modo que permitan diferenciar a este objeto de otros. Estas características, además, deberán ser comprensibles para los algoritmos o técnicas computacionales de minería de datos.

En el proceso de extracción de características se crea un vector por objeto (vector de características) o un conjunto de vectores para varios objetos (vectores de características) en los cuales las observaciones iniciales se condensan con el objetivo de reducir la información redundante y la dimensionalidad de los problemas. Esta reducción de dimensionalidad es particularmente importante cuando se dispone de una gran cantidad de datos, ya que en algunos casos el procesamiento de grandes cúmulos de información conlleva unos tiempos de cálculo prohibitivos.

Como ejemplo ilustrativo del proceso supóngase que en un estudio médico se intenta hallar los vectores de características que contengan la información más relevante para predecir un tipo de cáncer. El primer paso consistirá en seleccionar las variables medidas con mayor información y descartar las de menor. Así por ejemplo, se podrían descartar, por tener una relación mínima o nula con la patología a identificar:

- La altura del paciente.

- La presión sanguínea.
- La condición social.

Por otro lado, sí deberían considerarse, por ejemplo:

- Si el paciente es fumador habitual.
- El historial médico familiar.
- La exposición a radiaciones.

Nótese que en este ejemplo, de las 6 variables iniciales únicamente se conservan 3. La información inicial ya ha sido acotada (con una mínima pérdida de información relevante a la enfermedad a diagnosticar). El siguiente paso en el proceso de extracción de características podría consistir en recopilar, de las 3 variables consideradas, únicamente los datos familiares de tres generaciones, la exposición a radiaciones y únicamente considerarlo fumador si consume más de 3 cigarrillos por día.

El proceso de extracción de características es particularmente importante en el caso del estudio de fenomenología física en fusión nuclear, ya que los dispositivos generan una gran cantidad de datos en cada experimento. En algunas máquinas pueden adquirirse más de 1000 señales en cada descarga, conteniendo cada una de ellas una alta cantidad de muestras. En consecuencia, los tiempos de cálculo, utilizando estas bases de datos masivas, suelen ser extremadamente prolongados.

En esta Tesis, los vectores de características se crearán a partir de los datos derivados de descargas experimentales. En algunos casos, será una única magnitud, como la temperatura medida en un punto específico. En otros casos, el vector de características contendrá la información no de uno, sino de varias magnitudes. A los vectores creados se les aplicarán técnicas de minería de datos con cuatro objetivos generales:

- Búsquedas en bases de datos masivas.
- Creación de modelos físicos a través de los datos.
- Clasificación de eventos físicos.

- Predicción de fenomenología.

En los diferentes capítulos de la Tesis (en los cuales se explican las aplicaciones de las técnicas a problemas específicos de física de plasmas) se detallará el proceso de extracción de características implementado.

### **2.3 TÉCNICAS DE REDUCCIÓN DE DIMENSIONALIDAD Y VISUALIZACIÓN.**

#### **Reducción de dimensionalidad.**

El primer paso habitual en las técnicas de reducción de dimensionalidad consiste en la búsqueda de las propiedades que mejor describen a un grupo de objetos. Una vez encontradas esas propiedades, es posible descartar las que menos los representan, consiguiéndose así una disminución de dimensionalidad. Aunque la reducción de dimensionalidad siempre conlleva una cierta pérdida de información, las técnicas se desarrollan para que ésta sea la menor posible. Estas representaciones del conjunto de objetos originales suelen tener como principales objetivos minimizar los tiempos de cómputo y condensar la información relevante de un gran conjunto de datos. Sin ser su fin principal, pueden además facilitar la visualización de objetos de alta dimensionalidad: un vector de dimensión 50 no podrá ser pintado en un plano, pero sí es posible reducir su dimensionalidad a 3 y representarlo mediante un punto en un espacio tridimensional.

#### **Visualización.**

Por otro lado, existen técnicas específicamente desarrolladas con el fin de visualizar datos de alta dimensionalidad. En su objetivo final, entonces, es donde difieren de las anteriores: se centran en encontrar una representación visual fidedigna de los datos y no en la disminución de la dimensionalidad para disminuir los tiempos de cálculo. En general, las técnicas de visualización ayudan a identificar estructuras o agrupaciones de los objetos a través de las cuales se favorezca su interpretación. A su vez, pueden aplicárseles otras técnicas de minería de datos tales como las de clasificación, agrupamiento (“clustering”), estimación de densidad de probabilidad, etc.

En las siguientes sub-secciones se explican los métodos de reducción de dimensionalidad (análisis de componentes principales, la descomposición de valor

singular, el escalado multidimensional) y de visualización (el pseudo Grand Tour y los Mapas Topográficos Generativos) utilizados.

### 2.3.1 TÉCNICAS DE REDUCCIÓN DE DIMENSIONALIDAD.

#### 2.3.1.1 Análisis de Componentes Principales.

El Análisis de Componentes Principales [16] (ACP o PCA del inglés “Principal Component Analysis”) es un método estadístico concebido por Pearson a principios del siglo XX con el propósito de reducir grandes conjuntos de datos. Este método lineal se basa en calcular transformaciones ortogonales de las variables originales para conseguir un nuevo conjunto de variables no correlacionadas. A estas nuevas variables, que son el resultado de combinaciones lineales de las anteriores, se las denomina componentes principales. Las componentes finales del PCA quedan ordenadas de forma que cada una de ellas aporte la máxima contribución a la suma de las varianzas de las  $n$  variables. Así, la primera componente principal es la que resume de mejor forma la información original (i.e. contribuye mejor a explicar la varianza total). La segunda componente principal es la combinación lineal de las variables originales que mejor resume la varianza restante, y así sucesivamente hasta explicar el total de la varianza. La suma de varianzas de las  $n$  componentes principales es igual a la suma de varianzas de las  $n$  variables originales.

En general un conjunto reducido de factores contiene la mayor parte de la variabilidad total mientras que el resto de factores suelen contribuir comparativamente poco. Debido a esto, suele elegirse un limitado número de factores, reduciéndose significativamente la dimensionalidad del problema.

Existen dos formas básicas de aplicar PCA. Una de ellas es un método que utiliza la matriz de covarianzas. Ésta se usa cuando los datos son se encuentran en el mismo rango de variación y presentan valores medios similares. Este fue el método utilizado en esta Tesis doctoral para la extracción de características y reducción de dimensionalidad de magnitudes medidas en los dispositivos de fusión y por lo tanto será el único detallado (en la siguiente Sección).

La otra forma emplea la matriz de correlación. Esta es útil cuando los datos no son dimensionalmente homogéneos o cuando la magnitud de las variables aleatorias medidas fluctúa demasiado.

### 2.3.1.2 Criterios para decidir el número de componentes.

Para estimar la cantidad de componentes a retener, suelen utilizarse 2 métodos:

- Gráfico de autovalores: en los cuales se representa el porcentaje de variación explicada frente al número de la componente (Fig. 2.1). En el eje de las ordenadas se registra el porcentaje de varianza total explicada. En el eje de las abscisas se coloca el número del componente según su orden de importancia de acuerdo a la varianza explicada. Para el análisis de esta gráfica se buscan puntos en los cuales el cambio de la pendiente sea notable y la abscisa correspondiente a este punto indicará el número de componentes a retener.
- Promedio de autovalores: debe calcularse el promedio de todos los autovalores y eliminarse aquellos que se encuentren debajo de esta media. Si se utiliza la matriz de correlaciones esta media es igual a 1. Con este criterio se suelen retener menos componentes que con el anterior.

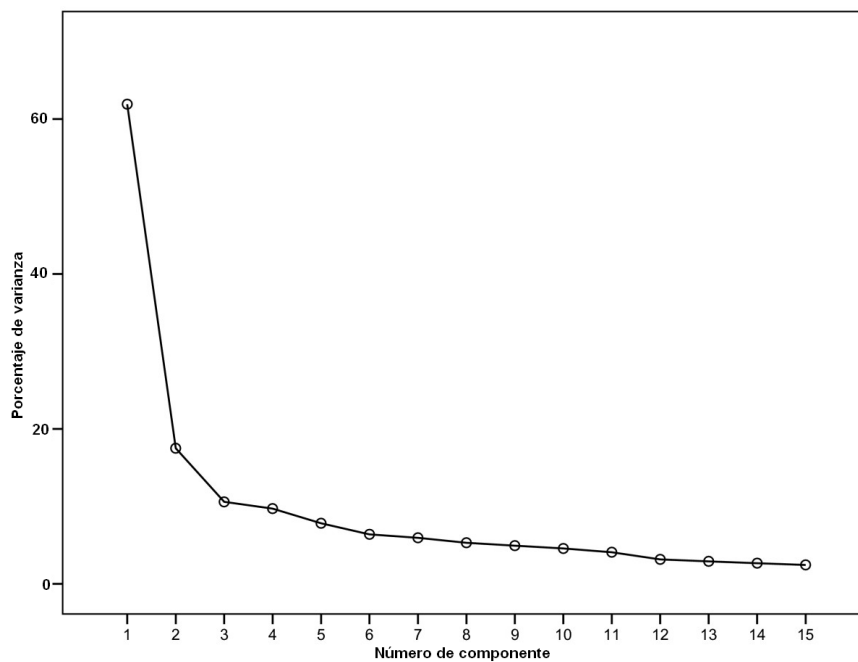


Fig. 2.1 Ejemplo de gráfico de autovalores.

*En general, las primeras componentes representan la mayor parte de la varianza. En este ejemplo podrían retenerse 2 o 3 componentes.*

### 2.3.1.3 Descomposición de valor singular (SVD).

La descomposición de valor singular (SVD, del inglés “Singular Value Decomposition”), es un método lineal de álgebra matricial vinculado al análisis de

componentes principales y suele obtener resultados similares. De hecho, proporciona una forma de encontrar las componentes principales sin tener que calcular explícitamente la matriz de covarianza [17].

Los datos son representados en una matriz  $X$ . La matriz  $X$  contendrá  $n$  filas y  $p$  columnas, donde  $n$  representa el número de objetos a analizar y  $p$  el número de características que describirán cada objeto.

El cálculo de la SVD de una matriz de datos  $X$  estará definida por:

$$X = UDV^T \quad (2.1)$$

Donde  $U$  es una matriz de dimensiones  $n \times n$ ,  $D$  es una matriz diagonal con  $n$  filas y  $p$  columnas, y  $V$  es una matriz de dimensiones  $p \times p$ . La matriz  $D$  contiene los valores singulares en su diagonal. Estos valores singulares son la raíz cuadrada de los autovalores de  $X^T X$ . Las columnas de la matriz  $U$  son llamados topos (ya que contienen las características estructurales de los objetos) o vectores singulares izquierdos (calculados como los autovectores de  $X^T X$ ). Similarmente las columnas de  $V$  son llamadas cronos (ya que suelen representar la evolución temporal de los correspondientes topos) o vectores singulares derechos [18].

Como con las componentes principales, los valores singulares se ordenan de mayor a menor, imponiéndose el mismo orden en las columnas de  $U$  y de  $V$ . La aproximación de menor orden de la matriz  $X$  original se denominará  $X_k$  y se obtendrá mediante:

$$X_k = U_k D_k V_k^T \quad (2.2)$$

Donde  $U_k$  es una matriz de dimensiones  $n \times k$  que contiene las primeras  $k$  columnas de  $U$ .  $V_k$  es la matriz de dimensiones  $p \times k$  cuyas columnas son las primeras  $k$  columnas de  $V$ .  $D_k$  es una matriz diagonal de dimensiones  $k \times k$  en la cual sus elementos diagonales son los  $k$  mayores valores singulares de  $X$ .

### 2.3.2 ESCALADO MULTIDIMENSIONAL (MDS).

El escalado multidimensional (más conocido por el acrónimo MDS, del inglés MultiDimensional Scaling) consiste en una serie de métodos no lineales que reducen la

dimensionalidad intentando preservar las distancias relativas de las observaciones cuando se transforman de un espacio de alta dimensionalidad a uno de menor dimensionalidad. El primer paso de MDS consiste en medir la proximidad de las observaciones para cuantificarlas. Es común entonces referirse en términos de similitud o disimilitud entre los objetos de entrada. Una medida de disimilitud entre dos objetos, llamémoslos  $r$  y  $s$ , puede denotarse como  $\delta_{rs}$  y la similitud entre ellos como  $s_{rs}$ . El método puede subdividirse en dos grupos [19]: los basados en métrica y los no métricos. Estos difieren principalmente en los procedimientos utilizados para la transformación de las disimilitudes  $\delta_{rs}$  en el espacio original a las disimilitudes en el espacio de dimensión reducida  $d_{rs}$ .

La técnica métrica utilizada en esta Tesis y detallada en la siguiente Sección se llama MDS clásico.

#### 2.3.2.1.1 MDS clásico.

Supongamos que las disimilitudes  $\delta_{rs}$  calculadas de los objetos de entrada  $p$ -dimensionales y las  $d_{rs}$  en un espacio de menores dimensiones pueden vincularse de la siguiente manera:

$$d_{rs} = f(\delta_{rs}) \quad (2.3)$$

Donde  $f$  es una función monótona continua. Para llegar a esta función puede utilizarse un método de optimización. Es necesario entonces plantear una función objetivo que cuantifique las discrepancias entre  $d_{rs}$  y  $f(\delta_{rs})$ :

$$\sqrt{\frac{\left( \sum_r \sum_s (f(\delta_{rs}) - d_{rs})^2 \right)}{\text{factor de escalado}}} \quad (2.4)$$

El factor de escalado más utilizado suele ser:  $\sum_r \sum_s d_{rs}^2$ .

En el MDS clásico las medidas de proximidad en el espacio original están basadas en la distancia Euclídea. Se puede demostrar que mediante estas distancias es



posible hallar una solución para una configuración en un espacio de menor dimensionalidad. La técnica puede resumirse en cinco pasos:

1- Llámese a la matriz que contiene las disimilitudes  $\delta_{rs}$   $\Delta$ . Debe encontrarse una matriz  $q$  donde cada uno de sus elementos sean dados por:

$$q_{rs} = -\frac{1}{2}\delta_{rs}^2 \quad (2.5)$$

2- Se calcula la matriz de centrado  $H$ :

$$H = I - n^{-1}OO^T \quad (2.6)$$

Donde  $I$  es la matriz identidad de dimensiones  $n \times n$  y  $O$  es un vector de  $n$  unos.

3- Encontrar la matriz  $B$  mediante:

$$B = HQH \quad (2.7)$$

4- Determinar los autovectores y autovalores de  $B$ :

$$B = ALA^T \quad (2.8)$$

5- Los vectores de características en el espacio de menor dimensionalidad estarán dadas por:

$$X = A_d L_d^{1/2} \quad (2.9)$$

Donde  $A_d$  contiene los autovectores correspondientes a los  $d$  más grandes autovalores y  $L_d^{1/2}$  contiene la raíz cuadrada de los  $d$  más grandes autovalores a lo largo de la diagonal. Si se elige una dimensionalidad  $d=2$  o  $d=3$ , los elementos pueden ser representados en gráficas bidimensionales o tridimensionales.

## 2.4 TÉCNICAS DE VISUALIZACIÓN.

### 2.4.1 TOURS VISUALES.

Los métodos basados en tours visuales permiten una visualización de los objetos multidimensionales en un espacio reducido a dos dimensiones. En particular el Pseudo Grand Tour (PsGT) [20], basado en el “torus grand tour” desarrollado por Asimov [21],

es un método computacionalmente eficiente que proporciona una secuencia de proyecciones de los objetos de alta dimensionalidad en un plano. Esta secuencia continua le otorga una apariencia de animación o vídeo. La secuencia de proyecciones puede detenerse cuando revele estructuras de interés. Si se desea, por ejemplo, diferenciar tres tipos de objetos, se buscará una proyección en la que se distingan tres grupos bien definidos. Por lo tanto, este tipo de técnicas debe ser, de algún modo, supervisada. Una forma de supervisión simple es mediante observación visual. Otra, es implementar en paralelo algoritmos que determinen automáticamente el paso en el que las proyecciones revelen estructuras o distribuciones de interés.

Para conseguirlo, PsGT genera 2 vectores,  $\alpha(k)$  y  $\beta(k)$ , donde  $k$  indica el número de iteración. Cada objeto podrá visualizarse desde todos los ángulos al ser multiplicado por los vectores creados (la diferencia del ángulo de proyección quedará determinada por el paso o iteración  $k$ ). El vector  $\alpha(k)$  proporcionará una de las coordenadas de la proyección del objeto en el plano para una iteración  $k$ . La otra se consigue mediante  $\beta(k)$ .

La secuencia de proyecciones se genera entonces multiplicando los objetos a representar por los vectores creados. En una primera iteración  $k$ , cada objeto (de alta dimensionalidad) es representado como un punto en un plano. En las siguientes, los objetos son representados desde otro ángulo, y así sucesivamente.

#### 2.4.2 MAPAS TOPOGRÁFICOS GENERATIVOS.

Aunque los Mapas Topográficos Generativos (conocidos con el acrónimo GTM, del inglés “Generative Topographic Maps”) son considerados como un método de aprendizaje, en esta Tesis han sido únicamente utilizados con fines de visualización y por lo tanto incluidos en esta Sección. GTM se basa en los mapas auto-organizados de Kohonen [22]. Los mapas de Kohonen son un tipo de red neuronal no supervisada, competitiva, distribuida de forma regular en una rejilla bidimensional en la que se representa la distribución de los datos de alta dimensionalidad. GTM, genera una distribución de los datos semejante [23] mejorando algunas limitaciones de los mapas de Kohonen. Estas limitaciones del método de Kohonen incluyen la falta tanto de pruebas matemáticas de convergencia como de una base teórica para la elección de los parámetros de aprendizaje y vecindad.

GTM, se define en términos de un “espacio latente” (bidimensional) en relación con el espacio de entrada (de alta dimensionalidad). Garantiza que la proximidad entre objetos en el espacio de alta dimensionalidad quede fiablemente representada por las distancias entre los puntos del espacio de baja dimensionalidad. Esto implica que objetos que se encuentren cerca en el espacio de entrada se transformarán en puntos, también cercanos, en el plano de representación. El plano de representación consiste en una rejilla cuyas dimensiones deben ser predefinidas por el usuario. Para calcular las transformaciones el método utiliza técnicas ampliamente reconocidas, tales como la estadística Bayesiana o el algoritmo “Expectation Maximization” [20].

El código utilizado en esta tesis fue desarrollado por Martínez [20] para su uso desde Matlab ®. Un ejemplo del resultado de una representación de objetos multidimensionales en un plano mediante GTM es mostrado en la Fig. 2.2.

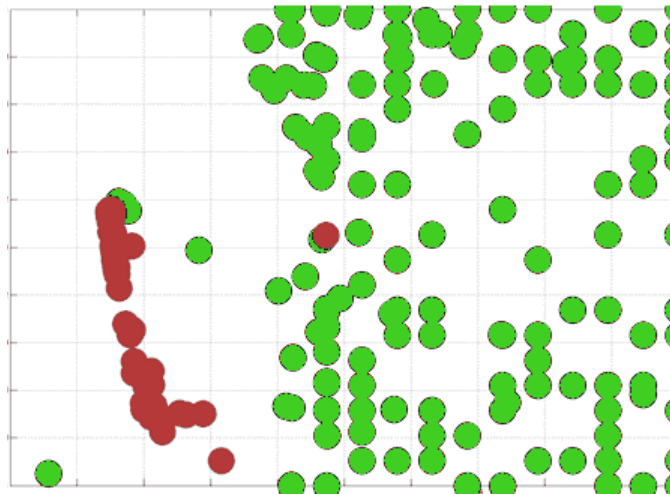


Fig. 2.2<sup>10</sup>. Ejemplo, en unidades arbitrarias, de una representación mediante GTM.

*Dos tipos de objetos de dimensión 13 son transformados en puntos en un plano.*

## 2.5 SISTEMAS DE APRENDIZAJE AUTOMÁTICO.

### 2.5.1 INTRODUCCIÓN.

El aprendizaje automático es una rama de la inteligencia artificial que se basa en la creación de programas capaces de generalizar comportamientos a partir de

información suministrada en forma de ejemplos. Su objetivo es automatizar algunas partes del método científico mediante técnicas matemáticas.

Los tipos de métodos de aprendizaje automático utilizados en esta Tesis son:

**Métodos de aprendizaje no supervisado:** se basan en el análisis, sin información a priori, de los datos de entrada. Los sistemas tratan tales datos como un conjunto de variables aleatorias y son capaces de determinar características comunes entre los objetos. El fin entonces es la clasificación o el agrupamiento de la información de entrada mediante medidas de similitud meramente matemáticas.

**Métodos de aprendizaje supervisado:** en ellos el algoritmo produce una función que establece una correspondencia entre las entradas y las salidas deseadas del sistema. Un ejemplo típico de este tipo de método es el problema de clasificación, donde el modelo construido estima las etiquetas de una serie de vectores de características utilizando una entre varias categorías o clases. La base de conocimiento con la que los sistemas construyen los modelos se forma con ejemplos de etiquetados anteriores.

En las próximas secciones se detallan los métodos no supervisados y supervisados utilizados en el desarrollo de esta Tesis.

## 2.5.2 MÉTODOS DE APRENDIZAJE NO SUPERVISADO.

### 2.5.2.1 Introducción.

Los métodos no supervisados estiman posibles relaciones entre los datos sin que sea necesario proveerles a priori información sobre ellos. Sin embargo, suele obtenerse mejores resultados cuando se aplica un preprocesado a ciertos problemas específicos a tratar, principalmente mediante la implementación de normalizaciones para que se asigne el peso deseado a las diferentes variables. Como resultado de los métodos no supervisados, se obtienen agrupamientos de la serie de objetos de entrada de forma que aquellos que pertenezcan a una misma clase presenten un alto grado de asociación.

Los métodos no supervisados pueden dividirse en dos grandes grupos:

- Los jerárquicos: son clasificadores que construyen un árbol en el cual cada nivel representa una partición de los datos. El nodo raíz del árbol corresponde a la partición menos refinada de los datos mientras que cada

partición posterior representa a un conjunto de objetos con características cada vez más similares. Cada nodo puede seguir subdividiéndose hasta que se considere necesario. Este tipo de métodos jerárquicos se conoce como divisivo. Si el árbol se construye en sentido inverso, desde las hojas a la raíz, mediante la fusión de clases, el método se denomina aglomerativo.

- Los no jerárquicos: en ellos el número de particiones en las que los datos deben agruparse se presupone como conocida. Sin embargo, no siempre la cantidad de clases en las que deben agruparse los datos se sabe *a priori*. En estos casos en los que no se conoce la cantidad de grupos en los que deben clasificarse los objetos de entrada suelen implementarse técnicas estadísticas en paralelo con las cuales se pueda estimar una distribución de los datos. Los métodos no jerárquicos se dividen en paramétricos y no paramétricos. Los métodos no paramétricos no asumen formas funcionales conocidas para la distribución de los datos, ya que es posible que las más usuales, como las Gaussianas, no se ajusten a las observadas.

En la siguiente Sección se detallará el método no jerárquico utilizado en el desarrollo de la Tesis: el algoritmo de K-medias (o K-means).

#### 2.5.2.2 K-medias o k-means.

El algoritmo de K-means fue desarrollado por MacQueen en 1967 [24, 25]. Este método agrupa y clasifica un conjunto de objetos de entrada en  $k$  clusters, donde la cantidad de clases  $k$  es fijada a priori. Dado el conjunto de observaciones iniciales (Fig. 2.3.a), el algoritmo iterativo se basa en asignar inicialmente  $k$  centroides (en el ejemplo  $k=3$ ), es decir un centroide por cada cluster (Fig. 2.3.b). Esta asignación inicial debe ser cuidadosa porque una ubicación distinta de los centroides puede producir diferentes resultados finales. Como norma general deben ser ubicados lo más separados unos de otros para obtener resultados satisfactorios. El siguiente paso consiste en asociar cada objeto al centroide más próximo (Fig. 2.3.c). Cuando esta asignación se ha completado para todos los objetos se concluye el primer paso. En este punto es necesario calcular nuevamente los centroides como baricentros de los clusters resultantes del paso anterior

(Fig. 2.3.d). Se obtienen como resultado nuevos  $k$  centroides cuya ubicación ha variado. Una vez más se asocian los objetos a los baricentros calculados, generándose así el bucle iterativo. Como resultado de cada iteración, los centroides modifican su posición anterior. El proceso termina cuando las posiciones de los centroides deja de variar, encontrando un equilibrio y agrupando los datos en  $k$  clusters (Fig. 2.3.f).

Para ello, el algoritmo utilizado por el método minimiza una función objetivo  $J$ , en este caso una función de error cuadrático que se define como:

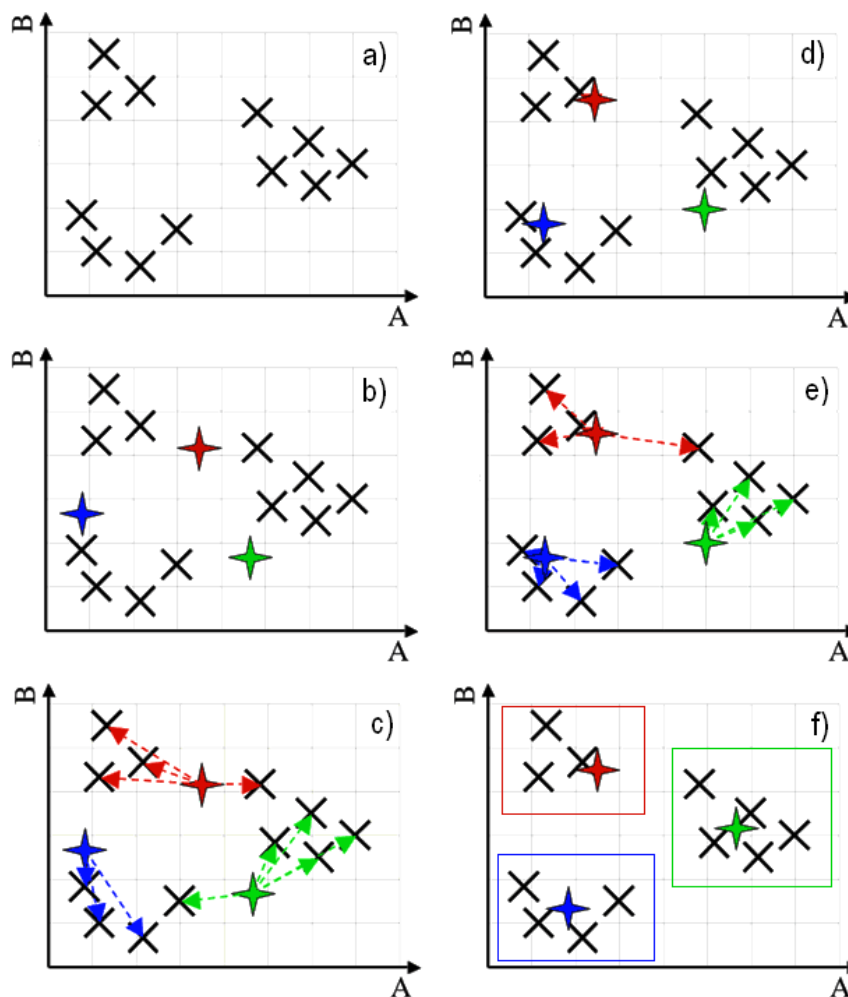


Fig. 2.3. Evolución del agrupamiento mediante  $K$ -means.

a) Grupo inicial de objetos a clasificar. b) Se asignan tres centroides. c) A cada objeto se le asigna la clase del centroide más próximo. d) Se calculan nuevamente los centroides como el baricentro de los objetos de cada clase. e) Se recalculan los centroides hasta la convergencia, en f), donde los baricentros no modifican su ubicación.

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_i\|^2 \quad (2.10)$$

Donde  $\|x_i^{(j)} - c_i\|^2$  es una medida de distancia elegida entre el objeto  $x_i^{(j)}$  y el centroide  $c_i$  a manera de indicador de distancias entre los  $n$  objetos y sus respectivos baricentros.

## 2.5.1 MÉTODOS DE APRENDIZAJE SUPERVISADO.

### 2.5.1.1 Introducción.

El aprendizaje supervisado puede ser utilizado para resolver problemas de regresión o de clasificación. En esta Tesis, los métodos son utilizados para clasificación. Los problemas de clasificación tienen como objetivo construir un modelo a partir de un conjunto de datos de entrada cuya clase se conoce. Una vez generado un modelo mediante ejemplos de objetos de entrada y la clase a la que pertenecen (entrenamiento), el sistema de clasificación debe ser capaz de actuar independientemente. Es decir, que el sistema creado debe ser capaz de estimar la clase de nuevos objetos (etapa de pruebas).

Los métodos consisten, entonces, en dos etapas. En la primera etapa de entrenamiento se introduce un grupo de objetos de entrada  $\bar{x}$  y el supervisor proporciona la clase y a la que pertenecen (ver Fig. 2.4). A partir de estos objetos-ejemplo y sus correspondientes clases, el modelo es generado.

En la segunda etapa, llamada fase de pruebas o de test, el sistema de clasificación debe ser capaz de estimar la clase de nuevos objetos (diferentes a los de entrenamiento pero provenientes de la misma fuente de información). Si el sistema ha sido entrenado debidamente, sus predicciones  $\hat{y}$  deben asemejarse a las que obtendría el supervisor. Para que el aprendizaje sea correcto, éste debe ser capaz de generalizar a partir de los ejemplos concretos de la etapa de entrenamiento, y por lo tanto en ella deben proporcionarse una cantidad de objetos suficientemente variada y cuantiosa. En la etapa de test también es necesario realizar pruebas con un número suficientemente amplio de ejemplos para verificar la fiabilidad de las clasificaciones.

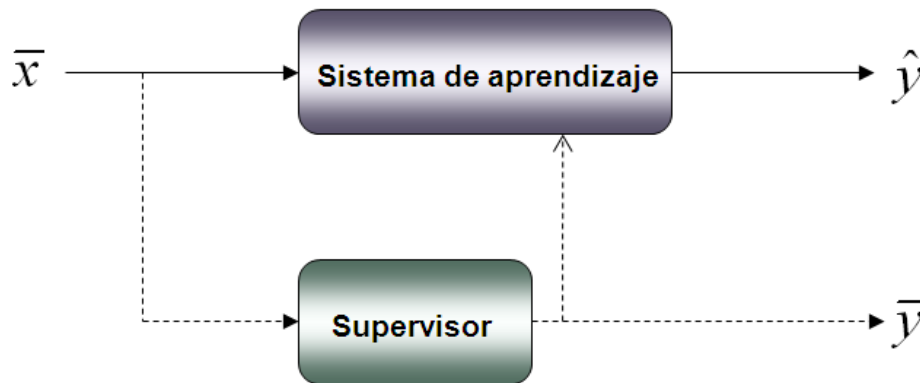


Fig. 2.4 Esquema de aprendizaje de un sistema supervisado.

Durante la etapa de entrenamiento un supervisor le indica al sistema de aprendizaje la clase a la que pertenece cada objeto de entrada. Una vez entrenado el sistema, este es capaz de distinguir automáticamente la clase a la que corresponden nuevos objetos.

En esta Sección se detallarán los dos métodos de aprendizaje no supervisado utilizados en el desarrollo de esta Tesis: los árboles de clasificación y regresión y las máquinas de vectores soporte.

#### 2.5.1.2 Árboles de clasificación y regresión.

Los árboles de clasificación constituyen uno de los métodos de aprendizaje inductivo supervisado no paramétrico más utilizados en la actualidad. Representan el conocimiento adquirido durante el aprendizaje mediante particiones recursivas de los datos. Estas particiones se traducen en una organización jerárquica que puede modelarse mediante una estructura arborescente. Cada nodo interior del árbol contiene una pregunta sobre un atributo concreto y cada nodo hoja se refiere a una decisión. La clasificación de patrones se realiza entonces mediante una serie de preguntas sobre los valores de sus atributos. Se comienza por el nodo raíz y se continúa el camino determinado por las respuestas a las preguntas de los nodos internos hasta llegar a un nodo hoja. La etiqueta asignada a esta hoja es la que se asignará al patrón a clasificar.

En particular el CART (sigla proveniente del inglés “Classification and Regresión Trees”) [26] se caracteriza por realizar particiones binarias y por utilizar una



estrategia de poda<sup>11</sup> basada en un criterio de coste-complejidad. Para ello, durante la fase de aprendizaje, el CART construye el árbol a partir de un conjunto de ejemplos de entrenamiento  $S$ . Este conjunto contendrá variables correspondientes a cada objeto y la clase respectiva. Para la construcción del árbol (Fig. 2.5), los algoritmos del CART inspeccionan todo el conjunto de objetos de entrenamiento con el objetivo de saber qué variable  $y$  y qué valor de dicha variable divide mejor los datos en las clases correspondientes. Cada pregunta sobre alguna variable de los objetos (realizadas en nodos interiores) dividen el conjunto de objetos en dos, de forma que cada nodo filial resultante maximice la “pureza” o minimice la “impureza” de su nodo parental. Para calcular la pureza existen diferentes fórmulas [26], pero la más aceptada y usada para el cálculo de la pureza en el nodo  $i$  es la siguiente:

$$I_G(i) = 1 - \sum_{j=1}^m p(i, j)^2 = \sum_{j \neq k} p(i, j) p(i, k) \quad (2.11)$$

donde  $m$  es el número de clases.  $p(i, j)$  representa la probabilidad de la clase  $j$  en el nodo  $i$  y  $p(i, k)$  la probabilidad de la clase  $k$  en el nodo  $i$ . La ecuación alcanza su mínimo (cero) cuando todos los objetos son clasificados en una única categoría.

La etapa de pruebas se realiza una vez creado el árbol. Un nueva serie de objetos responden a las preguntas asociadas a cada nodo hasta llegar a un nodo terminal. Allí, el nuevo objeto es clasificado como perteneciente a una de las clases predefinidas.

### 2.5.1.3 Máquinas de vectores soporte.

#### 2.5.1.3.1 Introducción.

Las máquinas de vectores soportes o SVM (siglas provenientes de “Support Vector Machines”) fueron desarrolladas por Vapnik en un principio para clasificación de datos y luego extendidas para su aplicación en problemas de regresión [27]. SVM se basa en el principio de minimización del riesgo estructural (SRM del inglés, “Structural Risk Minimization”) y los modelos que obtiene presentan la ventaja de depender únicamente de una pequeña proporción de los datos de entrenamiento, llamados

---

<sup>11</sup> El árbol creado mediante el CART tiene el nombre de “árbol máximo”. La poda consiste en la eliminación de las ramas más alejadas del nodo raíz y es necesaria para reducir el tamaño y la complejidad de árbol.

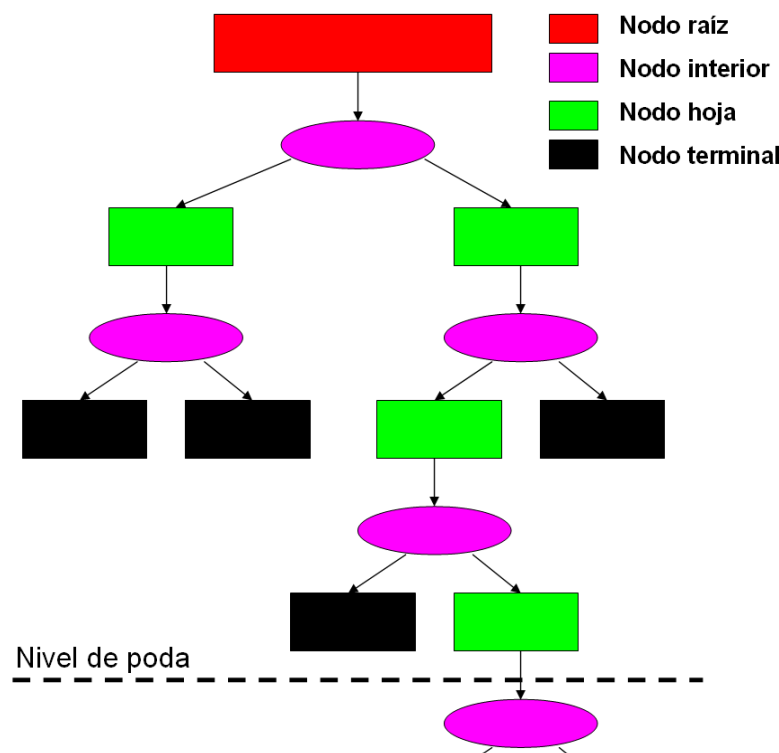


Fig. 2.5 Representación esquemática de un árbol de clasificación creado por el CART.

El árbol se divide en función de las condiciones impuestas por los nodos interiores (si una variable supera o no cierto valor). La subdivisión continúa hasta alcanzar nodos terminales o el nivel deseado de poda del árbol.

vectores de soporte, (SV por sus siglas en inglés, "Support Vectors"). Los clasificadores y modelos de regresión creados poseen una notable capacidad de generalización, necesitan pocos parámetros de ajuste y la estimación de los parámetros generales se realiza a través de la optimización de una función de coste convexa, lo cual evita mínimos locales. Al depender su solución únicamente de los SV, el modelo final puede ser escrito como una combinación de un número pequeño de vectores de entrenamiento.

### 2.5.1.3.1.1 Máquinas de vectores soporte lineales.

Para clasificación, SVM emplea los objetos de entrenamiento para construir, en un problema binario, un hiperplano<sup>12</sup> que los separa según la clase a la que pertenecen, con el objetivo de maximizar el margen (o la distancia) entre las clases.

Dado los objetos de entrenamiento  $x$  y sus correspondientes clases  $y$ :

<sup>12</sup> Se determina hiperplano a la generalización de un plano en un espacio de cualquier dimensión.

$$\{x_i, y_i\}; i = 1, \dots, l; y_i \in \{+1, -1\}; x_i \in \mathfrak{R}^d \quad (2.12)$$

La ecuación del hiperplano que divide las clases  $y_i$  quedará definida por:

$$\bar{w} \cdot \bar{x} - b = 0 \quad (2.13)$$

donde  $w$  es un vector normal al hiperplano de separación y  $b$  representa el término independiente.

Definiendo  $d_+$  y  $d_-$  como la distancia del hiperplano de separación al objeto más cercano de clase positiva y negativa respectivamente, se llamará margen de separación a la suma de las distancias  $d_+$  y  $d_-$  (Fig. 2.6). Estas distancias en valor absoluto, son idénticas:  $d_+ = d_- = \frac{|1|}{\|\bar{w}\|}$ , y por lo tanto el margen total a maximizar será

$$\frac{|2|}{\|\bar{w}\|}.$$

Además, en un problema linealmente separable, debe imponerse la restricción de que ningún objeto se encuentre entre los márgenes:

$$x_i \cdot \bar{w} - b \geq +1 \text{ para } y_i = +1 \quad (2.14)$$

$$x_i \cdot \bar{w} - b \leq -1 \text{ para } y_i = -1 \quad (2.15)$$

Las ecuaciones anteriores pueden reescribirse como:

$$y_i(x_i \cdot \bar{w} - b) - 1 \geq 0 \quad \forall i \quad (2.16)$$

La función objetivo deberá maximizar el margen con las restricciones previamente mencionadas. SVM, resuelve este problema de optimización cuadrática<sup>13</sup> mediante la implementación de los multiplicadores de Lagrange<sup>14</sup>. Como resultado se

<sup>13</sup> Tipo de optimización que utiliza una función cuadrática dependiente de variables con restricciones.

<sup>14</sup> Variables utilizadas en métodos de optimización que permiten hallar máximos o mínimos de una función objetivo con restricciones.

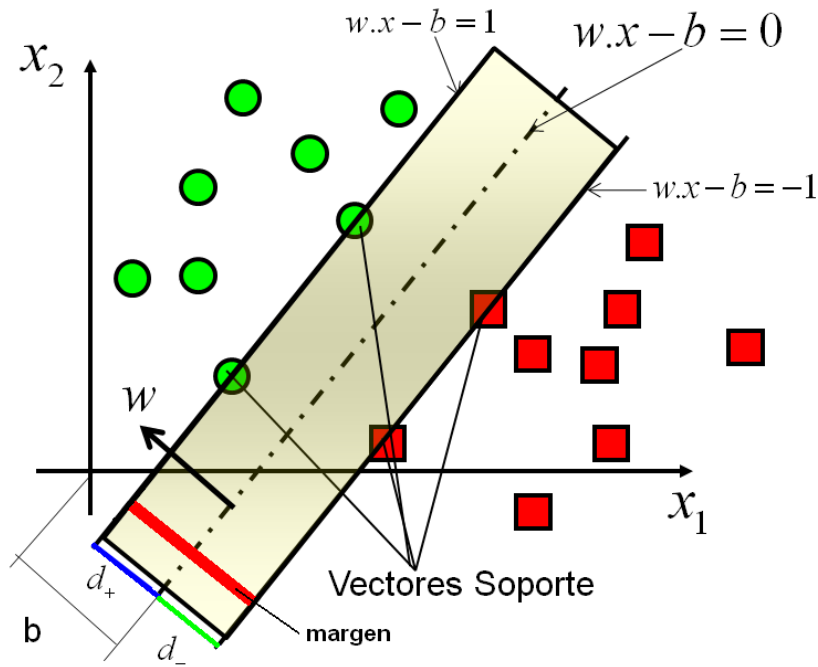


Fig. 2.6 Ejemplo de separación lineal con SVM.

El hiperplano divide las dos clases de objetos maximizando el margen de separación.

obtiene una función de decisión (la ecuación del hiperplano) la cual depende únicamente de los vectores soporte (Fig. 2.6).

Esta es:

$$D = \left( \sum_{i=1}^{N_{sv}} \alpha_i y_i x_i \right) - b \tag{2.17}$$

donde:

$$\bar{w} = \sum_{i=1}^{N_{sv}} \alpha_i y_i x_i \tag{2.18}$$

siendo  $N_{sv}$  el número de vectores soporte y los  $\alpha_i$  los multiplicadores de Lagrange. Estos multiplicadores, calculados en el proceso de optimización, serán iguales a 0 para todos los objetos que no sean vectores soporte.

### 2.5.1.3.1.2 Casos no linealmente separables.

Es posible flexibilizar el cálculo del hiperplano de separación, permitiendo que algunos de los objetos de entrenamiento se encuentren dentro del margen entre clases a maximizar. Para ello es necesaria la introducción de “slack variables”. Estas variables, durante el proceso de optimización, penalizan a los objetos situados dentro de los márgenes, por lo que se llega a una situación de compromiso entre el cálculo del hiperplano con mayor margen entre clases y los factores de penalización.

#### 2.5.1.3.2 Máquinas de vectores soporte no lineales.

La separación lineal no siempre es posible para todas las distribuciones de los objetos en el espacio de entrada [28]. La solución implementada por SVM consiste en llevar los datos a una dimensión mayor mediante el uso de funciones de transformación llamadas Kernels. En este nuevo espacio de mayor dimensionalidad, es posible calcular un hiperplano capaz de separar linealmente las clases. Así, la función de decisión quedará definida como:

$$D = \sum_{i=1}^m \alpha_i K(x_i, x) \quad (2.19)$$

siendo K la función Kernel de transformación.

Existen tantas funciones Kernel como se deseen, ya que pueden crearse *ad hoc*, aunque las más utilizadas suelen ser:

- Kernel polinómico:  $K(x_i, x_j) = (x_i \cdot x_j)^d$ , donde  $d$  indica la dimensión del polinomio.
- Kernel RBF (“Radial Basis Function”):  $K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$ , donde sigma es el parámetro que define la forma de la función radial.

Un ejemplo gráfico de cómo un problema no linealmente separable es transformado en uno separable (en un espacio de mayor dimensionalidad) es representado en la Fig. 2.7.

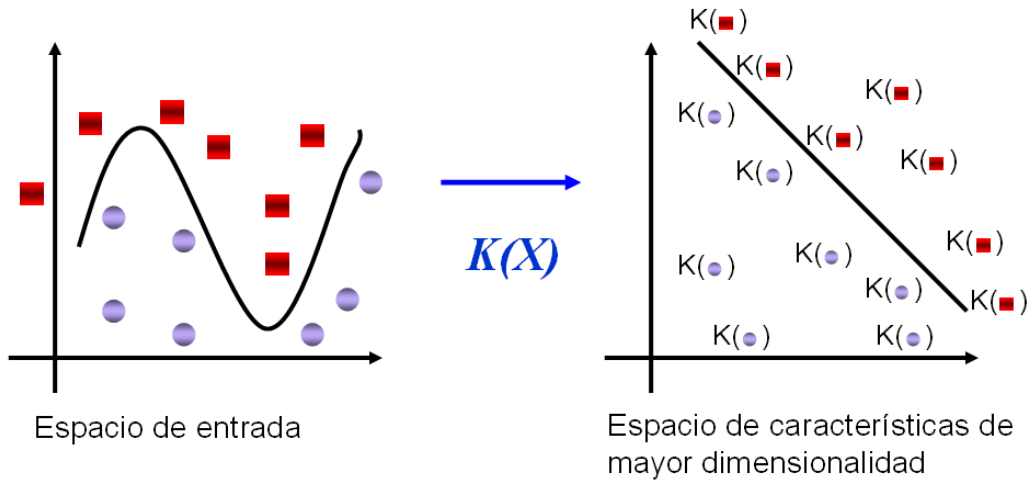


Fig. 2.7 Utilización de función Kernel.

Los datos, no separables linealmente en el espacio de entrada, son transformados a un espacio de mayor dimensionalidad donde puede encontrarse un hiperplano lineal que los divida según su clase.

---

## **3 Reconocimiento estructural de patrones en grandes bases de datos**

---

### **3.1 INTRODUCCIÓN Y ESTADO DEL ARTE.**

En varios campos de la ciencia moderna, desde la astrofísica a la ingeniería, la importancia de extraer y de interpretar debidamente la información almacenada en las bases de datos se ha incrementado continuamente. Esto se debe principalmente a que los avances impulsados por las nuevas tecnologías permiten almacenar extensas cantidades de datos. Estas nuevas capacidades han cambiado radicalmente la cantidad de información que puede ser guardada tras una medición. Sin embargo, no suelen incrementarse en esa proporción el uso de técnicas para el aprovechamiento de tal información. En un estudio de la Universidad de Berkeley se calculó que la producción de datos del año 2002 fue del orden de 4 exabytes [29]. Varios experimentos relacionados con investigaciones de fenómenos físicos siguen esta tendencia a almacenar significativas cantidades de información. Por ejemplo, el telescopio espacial Hubble ya ha enviado Tbytes de datos a la Tierra y se espera que el gran colisionador de hadrones, más conocido como “LHC” del CERN, produzca cantidades del orden de los 100 Pbytes.

Los plasmas generados en dispositivos por confinamiento magnético no son la excepción a esta tendencia. En particular, el Tokamak JET es capaz de producir pulsos de hasta 60 segundos de duración [15]. En cada uno de ellos la cantidad de datos puede

llegar a alcanzar medio Tbyte, cantidad que incluso puede incrementarse en las campañas experimentales planeadas para los próximos años [30]. Un variado tipo de señales (formas de onda, secuencias de imágenes, perfiles radiales, líneas de nivel) son adquiridas en tiempo real a través de los diagnósticos. Luego, en diferido, las señales son analizadas y si es necesario reprocesadas o validadas por expertos antes de guardarse en alguna base de datos especializada. Dado que todas las medidas de cada experimento pueden contener relevancia científica, todos los datos adquiridos son almacenados.

En los dispositivos para el estudio de la fusión, tanto Tokamaks como Stellarators, los parámetros medidos por los diagnósticos producen señales similares ante comportamientos físicos parecidos. Cuando la estructura morfológica particular de una forma de onda (o de una parte de ella) es de interés y se desea conocer en cuáles descargas existe un fenómeno semejante, éste debe buscarse manualmente, descarga a descarga, dentro de la amplísima base de datos acumulada durante varias décadas. Aún más, es posible que el interés resida no en el comportamiento particular de una señal sino en encontrar fenómenos físicos de mayor complejidad. Estos comportamientos pueden requerir el análisis no sólo de una forma de onda sino de la combinación de más de una señal de las adquiridas durante la ejecución de cada experimento. Como puede deducirse, tales procesos de búsqueda, realizados mediante la inspección visual de cada descarga, pueden demandar una enorme cantidad de tiempo.

Salta a la vista la necesidad de desarrollar métodos automáticos para la búsqueda de información de interés dentro de las extensas bases de datos de dispositivos de fusión nuclear. Esta necesidad se ve incrementada ante la duración prevista de las descargas (30 minutos) en futuras máquinas tales como ITER o Wendelstein 7X. Tales descargas producirán un volumen de datos significativamente mayor al actual.

En este capítulo se proponen soluciones para efectuar búsquedas en grandes bases de datos. Estas búsquedas estarán orientadas a encontrar de formas de onda completa, porciones de señales y fenómenos físicos.

En la Sección 3.1 se revisa el estado del arte, considerando los trabajos relacionados con búsquedas de formas de ondas producidas por diagnósticos de fusión.



En la Sección 3.2 se introducen las bases de las técnicas morfológicas (o estructurales), a través de las cuales es posible realizar búsquedas de formas de onda completa y de porciones de señales.

La Sección 3.3, se centra en la aplicación de técnicas estructurales en JET. Se describe el desarrollo, en primer lugar, de mecanismos de búsquedas morfológicas. Luego se demuestra cómo es posible adecuar tales técnicas estructurales, mediante la imposición de condiciones en las búsquedas, con el fin de reconocer fenómenos físicos específicos (en particular los cortes en canales de temperatura y las transiciones de modo L a modo H).

Finalmente en la Sección 3.4 se discuten los resultados generales.

## **3.2 RECONOCIMIENTO DE PATRONES EN GRANDES BASES DE DATOS.**

### **3.2.1 BÚSQUEDAS EN BASES DE DATOS DE FUSIÓN: ESTADO DEL ARTE.**

El acceso a datos puntuales en grandes bases de datos, como las de fusión, es imprescindible para el aprovechamiento de la información guardada. Para ello es necesario desarrollar mecanismos de búsqueda automáticos, veloces y efectivos. A pesar de ello, solamente un reducido número de trabajos han sido desarrollados en este área. En investigaciones previas se propusieron dos tipos de búsquedas automáticas. Un tipo tenía como objetivo encontrar formas de onda o señales completas, proporcionando como resultado una lista ordenada por semejanza. El otro tipo de búsqueda proponía un método para buscar formas morfológicas dentro de las señales. También, como resultado de tales búsquedas se proporcionaba una lista de los patrones encontrados ordenados mediante un criterio de similitud. En el desarrollo de esta Tesis, se propuso posteriormente un método alternativo, capaz de encontrar comportamientos físicos. Este método se basa en que, para algunos fenómenos determinados, la evolución de los parámetros del plasma presenta patrones morfológicos definidos. Así, estableciendo la búsqueda de una morfología específica, es posible encontrar eventos físicos estructuralmente semejantes.

En las siguientes secciones, luego de describir cada problema tratado, se detallará:

- La extracción de características, que consiste en el procesado de los objetos (creando así un vector de características) con dos propósitos principales:
  1. Traducir las características de los objetos bajo estudio a un lenguaje comprensible por un ordenador.
  2. Reducir la dimensionalidad de los objetos, condensando la información redundante y destacando la de mayor importancia.
- Las medidas de similitud empleadas, que determinarán el grado de semejanza de unos objetos con otros. Estas se aplican siempre que el sistema deba proveer una lista ordenada por la similitud del patrón y el de referencia. En las búsquedas de fenómenos físicos estas medidas no son implementadas ya que en estos casos se desean encontrar eventos y no es necesario ordenarlos por similitud.

### 3.2.2 BÚSQUEDAS DE FORMAS DE ONDA COMPLETAS.

Por ser los primeros en implementar técnicas de búsquedas automáticas en base de datos de fusión, cabe destacar tres aplicaciones que permiten seleccionar formas de onda completas [31, 32, 33], proporcionando como resultado una lista ordenada de las descargas con las señales similares.

Los métodos fueron implementados para señales de los Stellerators LHC (“Large Helical Device”) y TJ-II y posteriormente para las producidas por el Tokamak JET. Los tiempos de búsqueda para encontrar señales similares deben ser reducidos tanto como sea posible. Para ello, en estos trabajos los sistemas de clasificación se basan en acotar el espacio de búsqueda, de modo que la comparación de la señal de referencia se realice únicamente con las que *a priori* se consideren más probables y no con todas las almacenadas. Consecuentemente se evita el alto costo computacional de recorrer de principio a fin la base de datos para una comparación individual entre cada señal de entrada con todas las previamente almacenadas. La reducción del espacio de búsqueda se logra mediante estructuras jerárquicas (como las de la Fig. 3.1) y por lo tanto representables en forma de un árbol, de cuyo primer nodo raíz se despliegan ramas

que conducirán a nuevas hojas, nodos o “clusters”. Para realizar tales divisiones se siguen criterios de separación que consigan agrupar en cada nodo objetos con características comunes. Estos criterios de división suelen definirse específicamente para cada sistema a desarrollar e incluso pueden diferir en cada capa, tal como se demuestra en [30]. Es posible que algunos nodos del árbol de clasificación generado contengan una cantidad todavía demasiado amplia de señales. En ese caso es posible continuar subdividiéndolos.

La creación del sistema de clasificación se realiza sólo una vez. Posteriormente, al generarse nuevas señales estas pueden ser agregadas al sistema de clasificación mediante el mismo procedimiento. Al ejecutarse una búsqueda, los sistemas son capaces de proveer una lista ordenada por semejanza de las señales encontradas

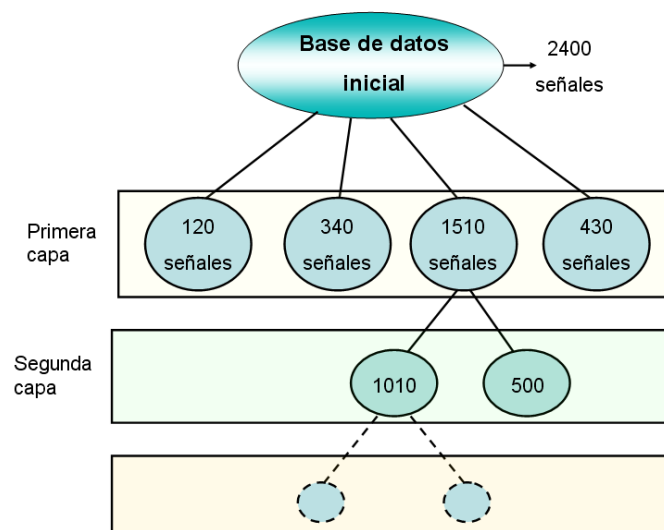


Fig. 3.1 Ejemplo de esquema de clasificación jerárquico.

*El sistema divide la colección de señales inicial (2400 en el ejemplo) en sub-grupos con características similares. Cada sub-grupo puede seguir subdividiéndose tanto como sea necesario. Las búsquedas se realizarán únicamente dentro del sub-grupo a priori más semejante a la señal de referencia.*

### 3.2.2.1 Extracción de características.

La extracción de características se realizó inicialmente mediante la transformada wavelet “Haar”. Esta puede ser calculada velozmente y además retiene tanto la información en frecuencia como la temporal. Esta extracción de características permite reducir la dimensionalidad del problema (de señales con decenas de miles de muestras a pocos coeficientes). Se decidió, en este caso, que 256 coeficientes eran suficientes para retener la información relevante. Como características morfológicas para los agrupamientos de se eligió la duración temporal de las señales. Agrupamientos basados en estas características han demostrado buenos resultados.

### 3.2.2.2 Medida de similitud.

La semejanza entre la señal de referencia y las señales encontradas se define mediante un criterio matemático de similitud entre las características de las señales a comparar. Esta medida puede considerarse como una distancia matemática. Dado que el ángulo entre dos vectores unitarios es una medida de la similitud entre ellos, en estos trabajos se eligió el producto escalar normalizado como método de comparación entre la señal de referencia (representada por su vector de características  $u_w$ ) y las almacenadas en la base de datos (con vectores de características  $v_w$ ). Una ventaja añadida de la medida elegida es que permite encontrar señales similares independientemente de amplitudes y polaridades. Este factor es importante en fusión, ya que en algunos casos las señales se adquieren con la polaridad invertida. En realidad, la medida de similitud es el valor absoluto de dicho producto escalar normalizado. Por tanto, dados dos vectores de características  $u_w$  y  $v_w$ , su similitud  $S$  queda definida por:

$$S_{u,v} = |\cos \alpha| = \frac{|u_w \cdot v_w|}{\|u_w\| \cdot \|v_w\|}, 0 \leq S_{uv} \leq 1 \quad (3.1)$$

Por lo tanto, cuanto más próximo sea  $S$  a 1, más semejantes se consideran las señales comparadas.

En este artículo las subdivisiones se proponen mediante agrupaciones obtenidas a partir de técnicas de análisis exploratorio de datos, en particular el método denominado Pseudo Grand Tour. En general, la cantidad de subdivisiones del sistema de clasificación dependerá tanto del tipo de información almacenada como del tamaño

de la base de datos en la que se deseen realizar las búsquedas. Una medida de similitud alternativa, basada en la transformada discreta de Fourier, se propone en 3.

### 3.2.3 INTRODUCCIÓN A LAS TÉCNICAS DE RECONOCIMIENTO MORFOLÓGICO PARA BÚSQUEDAS DE PORCIONES DE SEÑALES.

#### 3.2.3.1 Metodología general.

El método previamente descrito resulta adecuado para la búsqueda de formas de ondas similares a una de referencia. Sin embargo, es posible que el interés esté centrado únicamente en el comportamiento de una porción limitada de la señal y no en toda la onda. Para realizar eficazmente este tipo de búsquedas de porciones de señales, un trabajo pionero para señales de evolución temporal procedentes de diagnósticos del TJ-II fue desarrollado por Dormido-Canto y colegas [34] utilizando técnicas de reconocimiento morfológico. En ellas, se codifican con letras las señales de forma tal que queden representadas por un limitado número de características. La codificación resultante es almacenada en tablas. Para realizar las búsquedas, la porción de señal de referencia que desea encontrarse se codifica de la misma forma que se hizo con todas las señales de la base de datos. Las búsquedas entonces se reducen a encontrar caracteres: la cadena resultante se busca entre las almacenadas en la base de datos, obteniéndose una lista de todas las porciones semejantes ordenadas por similitud.

#### 3.2.3.2 Extracción de características.

El reconocimiento morfológico se basa en la idea de que las formas de las señales, por muy complejas que sean, pueden llegar a representarse mediante aproximaciones (lineales o de mayor orden) y aún así ser efectivas para el reconocimiento de formas particulares: comportamientos oscilatorios, dientes de sierra o incluso estructuras cualesquiera. Estos comportamientos pueden percibirse mediante una simple observación visual ya que a grandes rasgos representan la forma de una señal. Las características utilizadas en el reconocimiento morfológico y que posibilitan descomponer un problema en patrones más simples se llaman *primitivas*. Una traza de evolución temporal, por ejemplo, puede dividirse en segmentos lineales o en una sucesión de segmentos curvilíneos sin perder su morfología general. Una primitiva es el nombre asignado a estos segmentos. Cuanto mayor sea la cantidad de segmentos de aproximación, más parecida será la señal resultante a la onda original. Para aclarar el

concepto, supóngase que una serie temporal se desea aproximar por  $n$  segmentos lineales de duración  $k$  (aunque ajustes curvilíneos también pueden ser aplicados, tal como se propone en [36]). Un ajuste sencillo de la curva puede realizarse mediante los siguientes pasos:

*1- Calculando y predefiniendo la cantidad  $n$  de segmentos en los que se dividirá la señal (los criterios para esa determinación se discutirán posteriormente).*

*2- Efectuando un ajuste lineal de la señal original para cada uno de los segmentos  $k$ .*

*3- Dependiendo del ángulo resultante de cada aproximación como se indica en la Fig. 3.2, o de alguna otra propiedad como por ejemplo la diferencia de amplitud entre la primera y la última muestra de cada segmento  $k$ , se asignará una letra.*

Siguiendo estos pasos, la señal original, posiblemente compuesta por decenas de miles de muestras, quedará representada únicamente por  $n$  letras (Fig. 3.3). En este caso todos los segmentos tienen la misma longitud temporal  $k$ , por lo que llamaremos a esta forma de calcular las primitivas como técnica de Segmentos de Idéntica Duración (SID). Cada una de estas cadenas de letras que codifican las señales es almacenada en una base de datos. Esto puede considerarse como un inconveniente (ya que se requiere almacenar una cantidad extra de información, cuando las bases de datos son de por sí demasiado extensas). Sin embargo, al guardarse únicamente la codificación resultante, el espacio de almacenamiento extra requerido dependerá de la cantidad de caracteres con el que se codifiquen las señales. En la práctica no suele exceder el 2% de la base de datos original [35].

Un vez creada la tabla con las señales codificadas, las búsquedas de porciones de señales se convierte en una sencilla comparación de caracteres (entre los caracteres correspondientes a la porción de referencia y a los de la base de datos). Evidentemente, para este tipo de codificaciones no deben obviarse diversos factores que influirán en la

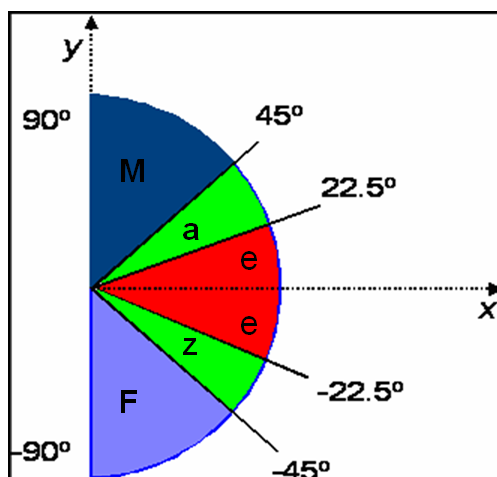


Fig. 3.2 Asignación de letras con respecto a las pendientes de los segmentos de ajuste.

efectividad de los resultados obtenidos, como por ejemplo la variedad de letras a asignar por cada segmento o la cantidad óptima de segmentos de codificación. Determinar automáticamente un óptimo es un problema complejo, porque un aumento de precisión en la búsquedas (codificando las señales con primitivas de menor duración, es decir con un mayor número de letras por cada señal) disminuirá la cantidad de formas similares (patrones) hallados, ya que estos contendrán una cadena de caracteres más larga que deberá coincidir con la de referencia. También influirán las características propias de las señales: si interesa encontrar eventos de variación rápida se necesitarán primitivas lo suficientemente cortas como para detectarlas. Normalmente la situación de compromiso está dada por la precisión requerida y la cantidad de señales similares que se encontrarán al realizar una búsqueda. En algunos casos se requieren sólo dos primitivas para obtener buenos resultados [37].

Alternativamente es posible, en el caso de conocerse o tener una estimación de la precisión requerida, determinar la duración de las primitivas en base al error de ajuste. A este tipo de primitivas se les puede llamar “Segmentos de Duración Variable” (SDV), aunque también ha sido nombrada de diversa forma (como “Constant Length Primitives” o Adaptive Length Primitives) en otras publicaciones [35, 36].

Para calcular este tipo de primitivas (Fig. 3.4) debe definirse un error máximo de ajuste ( $E_{max}$ ) permitido para cada aproximación lineal. La aproximación comienza con las primeras

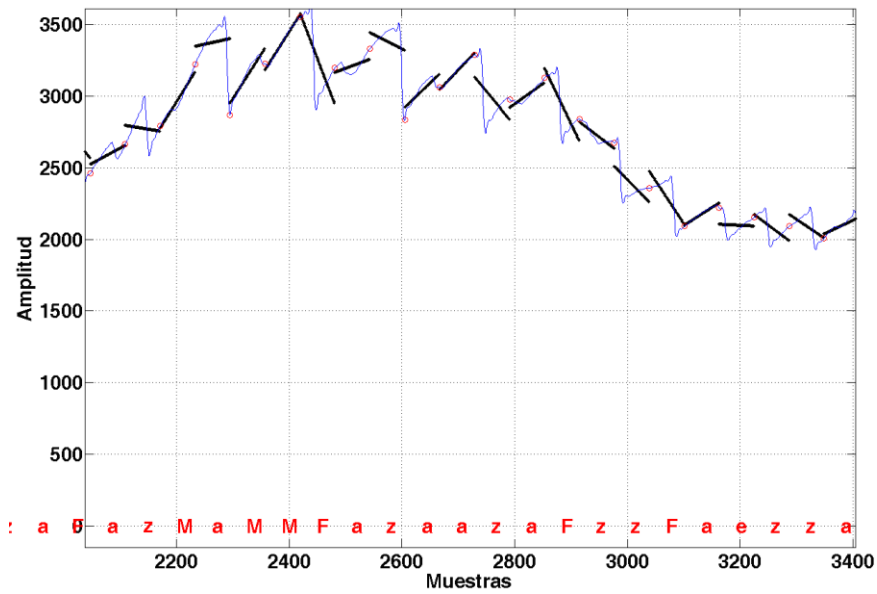


Fig. 3.3 Codificación “SID” de una señal. Todas las primitivas tienen una duración temporal idéntica.

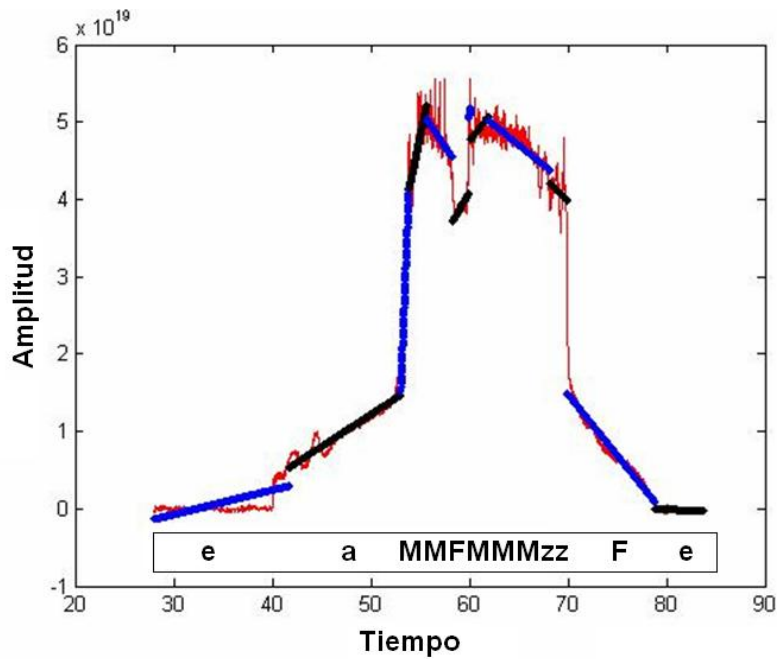


Fig. 3.4 Codificación “SDV” de una señal ejemplo.

La duración temporal de cada primitiva es variable y depende del error de ajuste.

muestras  $x_i$  y  $x_{i+t}$  (donde  $t$  es una cantidad constante predefinida de muestras de la señal original). Los pasos para la codificación se describen en el siguiente pseudo código:

MIENTRAS ( $i+t$  sea menor que el número total de muestras de la señal original)



1- Se calcula el ajuste lineal  $h(x)$  con las con las muestras comprendidas entre  $i$  y  $t$  de la señal original.

2- Se evalúa la diferencia (error) entre el segmento de aproximación y la función original.

SI (el error es menor que  $E_{max}$ )

Se incrementa  $t$  y se vuelve al paso 1.

FIN SI

SI (el error supera  $E_{max}$ )

Se da por concluida la primitiva actual.

$i=i+t+1$

Se vuelve al paso 1.

FIN SI

FIN MIENTRAS

Los resultados de una búsqueda en señales de temperatura se muestra en la Fig. 3.5. Esta señal describe la evolución de un canal de temperatura en el tiempo de duración de un experimento de TJ-II. Allí puede notarse que la estructura de las señales encontradas cumplen un patrón morfológico común aunque la duración de cada primitiva varíe en su longitud temporal. Para determinar la duración óptima de las primitivas se propuso en [34] una función de coste; esta función puede entenderse como un esfuerzo prometedor aunque quizás no definitivo para la solución de esta situación de compromiso.

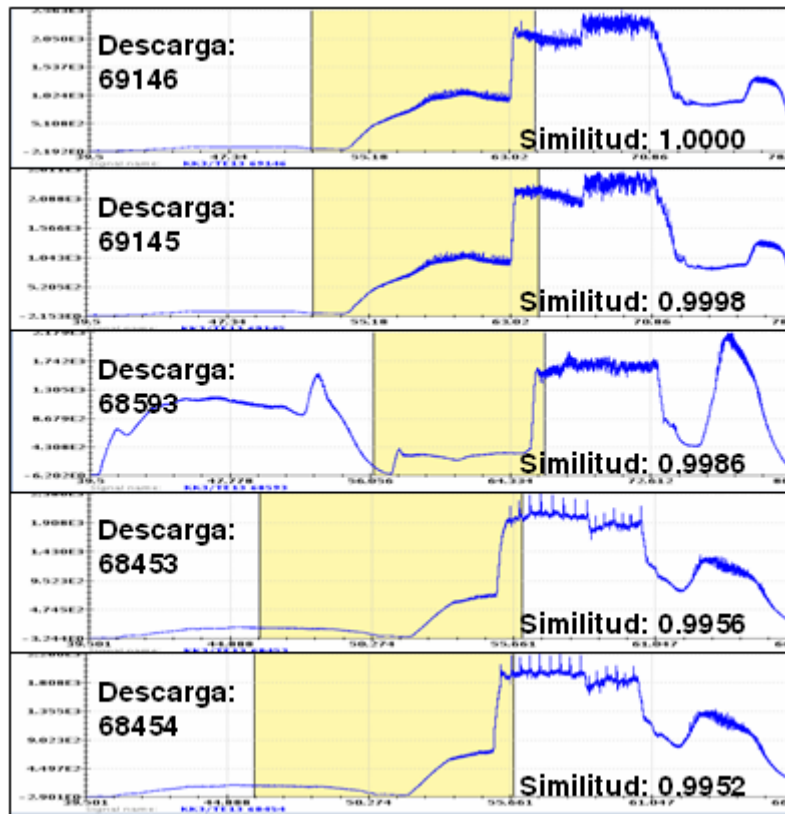


Fig. 3.5 Resultado de una búsqueda con primitivas SVD.

*La morfología en todos los casos es la misma: una caída, una pendiente positiva, una zona plana, un aumento abrupto y un final plano.*

### 3.2.3.3 Medida de similitud.

Las técnicas explicadas anteriormente codifican las señales reduciendo significativamente el espacio de almacenamiento requerido para describir el comportamiento de cada una de ellas. Una vez realizado el proceso, en la base de datos únicamente se almacena el resultado de la codificación: cadenas de caracteres y el número de descarga a la que pertenecen.

Las búsquedas entonces se realizan en tres pasos:

- El patrón de referencia buscado es seleccionado por el usuario. Automáticamente el sistema codifica el patrón en una cadena de caracteres.
- La cadena de caracteres se busca dentro de la base de datos.
- Una vez encontradas otras cadenas, los resultados de las búsquedas se ordenan según un criterio matemático de similitud: la distancia Euclídea.

Tal como se demuestra en [38, 39], con el mismo principio, estas técnicas pueden ser aplicadas al reconocimiento de patrones no sólo a señales de evolución temporal sino también a imágenes. En este caso, cada imagen se interpreta como una matriz de valores en la que cada vector fila se codifica como si se tratara de una señal de evolución temporal (ya que en definitiva la intensidad de los píxeles puede ser traducida a un valor de magnitud). Las búsquedas finales se realizan de manera semejante, aunque en este caso en vez de cadenas codificadas, se encuentran matrices de letras semejantes a la de referencia.

### **3.3 APLICACIONES DE TÉCNICAS ESTRUCTURALES EN JET.**

#### **3.3.1 BÚSQUEDA GENERAL DE PATRONES.**

##### **3.3.1.1 Extracción de características.**

Para esta aplicación específica de búsqueda de patrones en señales de JET hubo que considerar los tiempos de muestreo: estos no eran los mismos para todas las señales adquiridas. Por lo tanto, el proceso de extracción de características comenzó con un preprocesado de las señales. Cada una de ellas fue interpolada con el fin que la diferencia de tiempos entre muestras consecutivas fuese siempre la misma para todas las formas de onda, guardándose una cantidad limitada de muestras por señal. Como se explicó anteriormente, el reconocimiento estructural es útil para búsquedas morfológicas y la interpolación no afecta significativamente a la forma de las señales.

El radiómetro de JET consta de 96 canales para medir temperaturas. Cada uno de los canales genera una señal de aproximadamente 52000 muestras. Estas muestras se redujeron, mediante la interpolación, a 4096 y se codificaron con primitivas SID. Las primitivas se almacenaron (una por cada cuarto de segundo de señal) en una base de datos de postGre, la cual permite realizar búsquedas condicionadas de forma rápida y sencilla. La base de datos se pobló con cadenas de caracteres, donde cada letra representaba una primitiva.

El protocolo de comunicación implementado un TCP/IP Cliente/Servidor (permitiendo acceso concurrente y simultáneo de varios usuarios) usando Berkeley Sockets como API (“Application Program Interface”), como se puede observar en la Fig. 3.6. Para la selección de la porción de interés de la señal de referencia se desarrolló una interfaz gráfica (Fig. 3.7) accesible desde los JACs (Jet Analysis Clusters) de JET.

Una vez ejecutada la búsqueda, el software proporciona como salida una lista en la que se especifica el número de descarga y los tiempos en los que se encuentran dentro de cada señal hallada las estructuras semejantes a la seleccionada en la señal de referencia (Fig. 3.7.b).

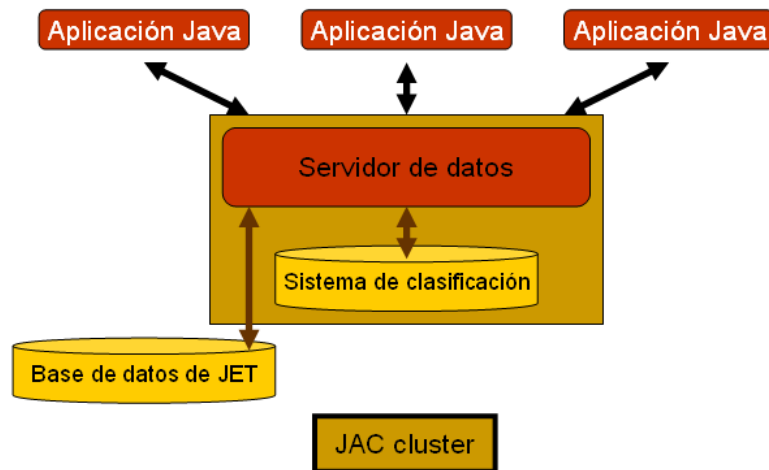


Fig. 3.6 Esquema del protocolo de comunicación empleado en JET.

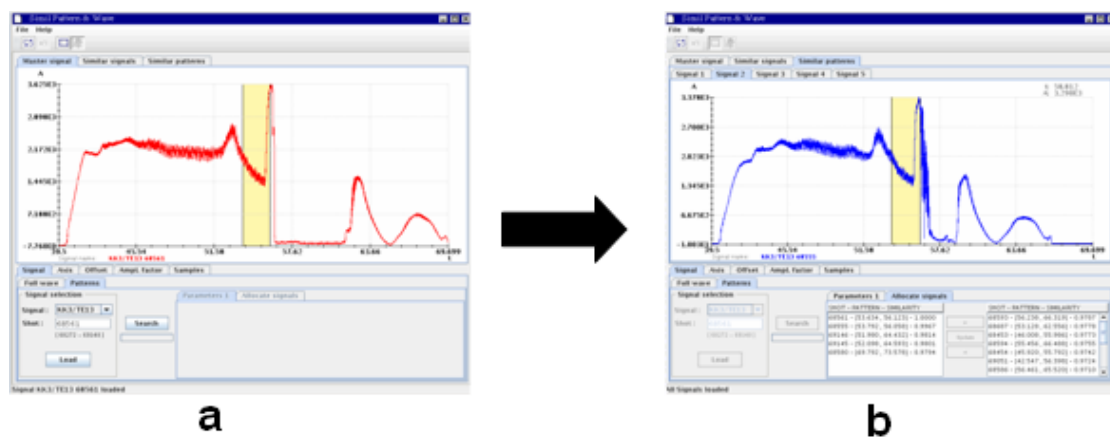


Fig. 3.7 Interfaz gráfica diseñada para la búsqueda de patrones dentro de señales.

a) Se selecciona la porción de interés (amarillo) de la señal de referencia. b) El sistema proporciona una lista ordenada de las porciones semejantes encontradas en la base de datos.

### 3.3.1.2 Medida de similitud.

La lista proporcionada está ordenada según la semejanza entre el segmento de señal de referencia y los encontrados. El criterio de similitud adoptado es la distancia euclídea entre ellos. Los tiempos requeridos para las búsquedas varían dependiendo en proporción al tamaño de la base de datos. Para el estudio detallado en [40] se requirieron 90 ms para una búsqueda en una base de datos de 501 señales.

## 3.3.2 BÚSQUEDA DE FENÓMENOS FÍSICOS ESPECÍFICOS.

### 3.3.2.1 Introducción.

El reconocimiento estructural juega un papel central en la distinción de comportamientos físicos cuando se analizan señales. Generalmente, una componente sinusoidal, una variación rápida o un pico en una forma de onda son suficientes como para identificar un fenómeno físico específico. En otras ocasiones, sin embargo, para detectar un comportamiento particular del plasma es necesaria la coincidencia de múltiples eventos en más de una señal.

Una herramienta rápida, efectiva, de fácil utilización, capaz de realizar búsquedas de fenómenos físicos puede ser de gran ayuda en el ámbito de la fusión. Por ello, y como prueba de principios, se desarrollaron en esta Tesis métodos morfológicos con el fin de buscar dos fenomenologías específicas.

### 3.3.2.2 Cortes en canales de ECE.

En el primer caso se utilizaron para la detección de los tiempos en los que ciertos canales de temperatura entran en corte. Determinar automáticamente estos tiempos, como se verá posteriormente, permite calcular otros parámetros, como la densidades a las que se producen. Esta aplicación sirve de ejemplo de reconocimiento de un fenómeno específico mediante la identificación de un único patrón. Las señales provienen de uno de los diagnósticos utilizados para medir la temperatura en plasmas de JET: el radiómetro heterodino [41]. Su principio de funcionamiento está basado en la emisión ECE. La antena del radiómetro realiza mediciones a determinadas frecuencias fijas (96 canales), lo que significa que cada uno de sus canales capta la emisión de un punto concreto del plasma. Una de las limitaciones de cualquier sistema ECE en plasmas de alta densidad es la aparición de cortes (reflexión interna de la radiación emitida), ya que a determinados niveles de densidad la emisión no llega a ser recibida

por la antena sintonizada a esa frecuencia. Esto se debe a que la propagación de las ondas en el plasma tiene una frecuencia de corte que depende de la densidad. Cuando la densidad entre la antena de ECE (ubicada en el lado de bajo campo del plasma) y la región de emisión supera los valores de corte, la radiación no se propaga hasta alcanzar el sistema receptor. Consecuentemente se produce una reducción significativa de las lecturas en los canales centrales, manteniendo un nivel de señal estacionario y casi nulo (canal en corte). Posteriormente, cuando la densidad cae bajo los valores de corte, la señal rápidamente vuelve a alcanzar su nivel adecuado. El tiempo en el que los canales no reciben emisión alguna puede variar en cada descarga y por lo tanto las “formas” a buscar no son siempre las mismas.

En términos simples, el comportamiento puede resumirse en tres etapas (Fig. 3.8):

1. Una caída brusca en la temperatura.
2. Un tiempo, variable, en el que la temperatura se mantiene muy baja y no varía significativamente en su amplitud.
3. Un incremento considerable de la temperatura al salir los canales del corte.

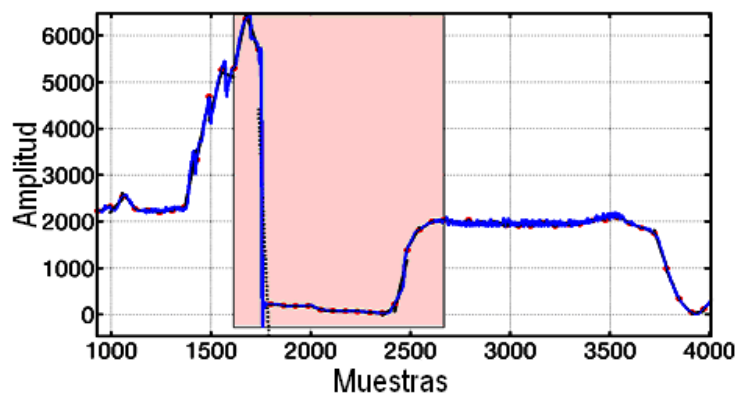


Fig. 3.8 Señal de temperatura, correspondiente al canal 13, que entra y sale de un corte.

#### 3.3.2.2.1 Extracción de características.

Para detectar este fenómeno se realizaron modificaciones a las búsquedas morfológicas, añadiendo a las técnicas de reconocimiento estructural el conocimiento sobre el comportamiento del plasma y su consecuencia en este diagnóstico. Teniendo en cuenta que en cada pulso se generan 96 señales de temperatura (una por cada canal del

radiómetro) y que cada una de ellas cuenta con 51200 muestras, fue necesario reducir la cantidad de datos a analizar para que el tiempo de las búsquedas fuera abordable. Para conseguirlo se obviaron los canales que medían en radios alejados del centro del plasma, ya que los cortes ocurren a altas densidades y por lo tanto los canales correspondientes a regiones próximas a la pared del dispositivo no son susceptibles a tales fenómenos. Para conocer exactamente la configuración de los canales en cada disparo se accedió a las señales procesadas `ppf\kk3\GEN` y `ppf\kk3\CPRF`. La primera contiene información general sobre el radiómetro. La segunda especifica el radio al que está midiendo la temperatura cada uno de sus 96 canales. En estas señales se buscaron los canales que medían en la regiones de mayor interés, dos segundos después que existiera calentamiento NBI. Con ese criterio se seleccionaron para cada descarga 2 señales de las 96. Estas dos señales representan mediciones en los radios del dispositivo que es más probable<sup>15</sup> que se produzcan cortes debido a modo ordinario o extraordinario de reflexión<sup>16</sup> respectivamente (el modo de medición queda especificado por la señal `ppf\kk3\GEN`).

Se analizó un total de 15850 señales de JET (7925 de cada uno de los dos canales de temperatura seleccionados, desde la descarga 60131 hasta la 69000). Para cada descarga, la pareja de señales fue procesada y codificada con la técnica SID (ver Sección 3.1).

La estructura que se desea encontrar no responde a una forma fija: los canales pueden entrar y salir del corte durante diferentes intervalos de tiempo. Por lo tanto fue necesario imponer condiciones en la búsqueda. La asignación de letras para la codificación fue la misma representada en la Fig. 3.2.

La base de datos en las que se almacenaron las cadenas fue, en este caso, MS Access, de fácil implementación, accesible desde Matlab y en la que se pueden imponer

---

<sup>15</sup> Según los expertos en el diagnóstico, aproximadamente a 3.1 metros para modo O y 3,5 metros para modo X.

<sup>16</sup> Modo O (ordinario): Modo de propagación de las ondas electromagnéticas en un plasma en dirección perpendicular al campo magnético externo, en el que el campo eléctrico de la onda es paralelo a dicho campo magnético.

Modo X (extraordinario): Modo de propagación de las ondas electromagnéticas en un plasma en dirección perpendicular al campo magnético externo, en el que el campo eléctrico de la onda es perpendicular a dicho campo magnético.

condiciones de búsqueda. Además de las letras que conforman cada cadena, se agregaron campos con los tiempos correspondientes a la ubicación temporal de cada primitiva almacenada. Estos campos, como se explicará posteriormente, son utilizados para determinar los instantes en los que los canales entran y salen del corte.

Las condiciones de búsqueda, basadas en las descritas anteriormente y que a grandes rasgos definen el fenómeno, fueron:

1. Una pronunciada pendiente negativa, es decir que la cadena debe iniciarse con la letra F (Fig. 3.2) que representa la mayor derivada negativa.
2. Un tiempo variable en el cual la pendiente sea prácticamente nula (es decir una cadena de longitud indefinida donde las letras representen derivadas pequeñas).
3. Una pronunciada pendiente positiva (letra M o a, las de mayores pendientes positivas).

#### 3.3.2.2.2 Resultados.

Tales comportamientos, sin embargo, no describen el fenómeno unívocamente, ya que cambios abruptos en la temperatura del plasma pueden ser malinterpretados como cortes en los canales de temperatura. No obstante, estos casos no suelen ser muy comunes. El porcentaje de canales en corte identificado correctamente fue del 92% de las descargas en las que se evidencia dicho fenómeno. Estos resultados son considerablemente satisfactorios si se tiene en cuenta la simplicidad del criterio de búsqueda implementado.

Un ejemplo de búsqueda puede observarse en la Fig. 3.9. Allí, se muestran dos resultados: uno correcto y uno erróneo. Nótese que la equivocación corresponde, al menos en su morfología, al patrón buscado.

Como además de las primitivas se guardaron los tiempos correspondientes a cada una de ellas, fue posible disponer de los instantes en los que los canales entran y salen del corte. Esta información fue utilizada para crear la Tabla 3.1. En ella puede observarse un ejemplo de identificación para un caso de modo ordinario y otro para



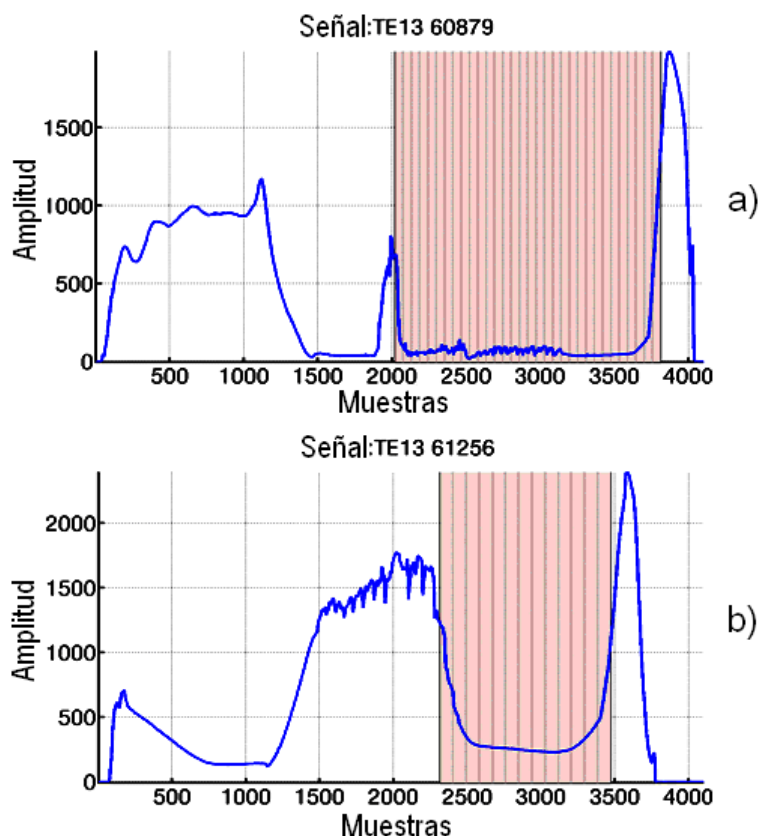


Fig. 3.9 Dos ejemplos de los resultados de las búsquedas de cortes en canales de temperatura.  
 a) Ejemplo de uno de los resultados correctos proporcionados por la búsqueda del fenómeno. b) Ejemplo de un error. Nótese en este caso que la estructura general se asemeja a la de los canales en corte.

Nº Canal	Nº Descarga	Modo	Instantes de Corte [s]		Densidades medidas [ $m^{-3}$ ]		Densidad Estimada [ $m^{-3}$ ]
			Inicio	Fin	Inicio	Fin	
6	66258	O	56.42	64	6.06e+19	6.58e+19	6.333e+019
1	66442	X	50.69	57.26	2.61e+19	2.99e+19	2.937e+019

Tabla 3.1. Representación de dos descargas en corte encontradas, una de ellas en modo O y la otra en modo X.

Se especifican para cada caso los tiempos de iniciales y finales en los que los canales entraron en corte, las densidades medidas en cada uno de esos instantes y la densidad estimada.

modo extraordinario. Allí son guardados los instantes en los que los canales entran y salen del corte. Acudiendo a la señal ppf/lidr/ne fue posible conocer la densidad medida

en esos instantes. Finalmente, y sabiendo la frecuencia a la que cada canal del radiómetro está sintonizado, fue posible estimar la densidad en el instante de corte mediante las siguientes fórmulas:

Para el modo ordinario (O):

$$ne = \frac{wc^2 \cdot \varepsilon_0 \cdot me}{e} \quad (3.2)$$

$$ne = \frac{wc^2 \cdot \varepsilon_0 \cdot me}{e} \quad (3.3)$$

Para el modo extraordinario (X):

$$ne = \sqrt{\left( wc \pm \frac{wce}{2} \right) - \frac{wce^2}{4}} \cdot \frac{me \cdot \varepsilon_0}{e^2} \quad (3.4)$$

Donde:

$$wce = \frac{e \cdot B}{me} \quad (3.5)$$

Siendo:

$\varepsilon_0$ , la constante de permitividad del vacío.

$me$ , la masa del electrón.

$e$ , la carga del electrón.

$ne$ , la densidad electrónica.

$wce$ , la frecuencia angular electrónica.

$wc$ , la frecuencia angular de corte.

$wpe$ , la frecuencia angular electrónica del plasma.

También se incorporan a la tabla las densidades integradas para la línea central medidas en cada uno de esos instantes de interés (extraídas de la señal PPF/LIDR/NE de JET).

### 3.3.2.3 Transiciones L-H.

Mediante el extenso estudio de las propiedades del plasma en diferentes Tokamaks se han identificado a través de los años diferentes regímenes de confinamiento. Conocidos como modos de confinamiento, estos regímenes usualmente ofrecen características únicas en cuanto a los procesos físicos experimentados por el plasma. Los modos más observados y estudiados son el modo óhmico, el modo L, el modo H y las configuraciones que incluyen barreras internas de transporte, las cuales mejoran sustancialmente el confinamiento [42]. Estos regímenes han sido no sólo observados sino además reproducidos en un gran número de dispositivos.

Inicialmente, la energía aplicada al plasma estaba predominada por el calentamiento óhmico (modo óhmico). Aunque los tiempos de confinamiento en este tipo de operación aparentaban ser prometedores, se descubrió luego que el puro calentamiento óhmico era insuficiente para alcanzar los altos parámetros requeridos por un reactor termonuclear. Mientras la temperatura del plasma aumenta, su resistencia disminuye, reduciendo consecuentemente el nivel de calentamiento que puede ser obtenido. El uso de técnicas adicionales para el calentamiento (generalmente inyección de partículas neutras o microondas) incrementaron ostensiblemente las temperaturas alcanzadas. En detrimento, estas intervenciones externas suele crear inestabilidades en el plasma, lo que produce una reducción del tiempo de confinamiento. Este tipo de régimen con calentamiento adicional se denomina modo L (“Low confinement Mode”), o modo de bajo confinamiento.

En año 1982 se descubrió en ASDEX que por encima de un cierto umbral de potencia de calentamiento, puede evidenciarse un alto gradiente de presión en el borde (denominado pedestal), un incremento en el tiempo de confinamiento de la energía de más de un factor 2 y la aparición de modos localizados en el borde, denominados ELMs (Edge Localized Modes) [43]. A este régimen se le llamó modo de alto confinamiento o modo H (del inglés “High Confinement”).

La determinación del instante en el que el plasma cambia de un régimen de bajo a alto confinamiento tiene un marcado interés científico y por ello se han desarrollado cuantiosos sistemas para su detección automática. Estos estudios están mayormente basados en métodos estadísticos o técnicas no lineales de aprendizaje, como “Decision and Regression Trees” (CART), lógica difusa, “Redes Neuronales Artificiales” (ANN), “Support Vector Machines” (SVM) o combinaciones de ellos. Estos métodos obtienen buenos resultados pero necesitan de una base de datos para su entrenamiento. El reconocimiento estructural de patrones puede ser usado como una aproximación sencilla al problema para identificar y acotar automáticamente el intervalo en el que tiene lugar la transición.

#### 3.3.2.3.1 Extracción de características.

La transición L-H (“Low mode to High mode transition”) suele identificarse visualmente en la mayoría de los casos mediante una inspección minuciosa de la señal DALPHA<sup>17</sup>, en particular la línea de visión del exterior desde el divertor de JET. La señal suele evidenciar la transición L-H mediante una caída en su amplitud en un periodo de tiempo de pocos milisegundos. Está compuesta por más de 130000 muestras, se encuentra afectada por ruido y sus magnitudes suelen variar drásticamente en cada experimento (las diferencias pueden ser de varios órdenes) o incluso en un mismo pulso debido a los ELMs. Sin embargo, esta rápida disminución de la magnitud de la señal no siempre es notoria y en muchos casos puede confundirse con el mismo ruido que afecta la señal. Si únicamente se analiza esta medición la correcta determinación del instante de la transición no es posible. Inevitablemente es necesaria otra señal o grupo de señales para añadir información que valide los resultados y ayude a alcanzar niveles de precisión razonables en la identificación de las transiciones.

En la Fig. 3.10.a se representa la codificación SDV aplicada a una señal DALPHA de JET. Puede notarse a simple vista que cuando ocurre la transición (aproximadamente a los 57.08 segundos), se produce una caída rápida y la longitud de la recta de ajuste es notablemente superior a la promedio. El criterio de búsqueda, se definió asumiendo que una transición L-H se detectará cuando la primera primitiva que

---

<sup>17</sup> Señal que mide la radiación con una longitud de onda de aprox. 650 nm. Esta radiación es producida cuando el electrón del átomo de deuterio experimenta una transición entre los niveles  $n=3$  y  $n=2$ .

represente una marcada pendiente negativa y una longitud de la recta de ajuste al menos dos veces superior a la longitud promedio. Para tener disponibles las longitudes de las rectas de ajuste, se incluyó un campo en la base de datos MS Access. Este campo especifica, además de la letra que describe cada primitiva, su longitud. Como se mencionó anteriormente, esta condición es necesaria pero no suficiente. Para disminuir los errores en la identificación se añadió el análisis detallado de la señal de densidad integrada de línea (GS/SS-DENS1 de JET), buscando dentro de ella pendientes positivas. Esto se debe a que un incremento en la densidad suele ser una consecuencia de la transición: al mejorarse el confinamiento menos partículas se pierden y la densidad cerca del centro del plasma aumenta. Se tuvo también en cuenta que los aumentos en la densidad no corresponden únicamente a un cambio de régimen, sino que pueden ser consecuencia de otros factores (inyección de gas, de partículas neutras o de pellets<sup>18</sup>, así como de un aumento de impurezas). Por lo tanto fue crucial la elección apropiada de las pendientes observadas en esta señal que evidencian una mejora del confinamiento y no otros posibles factores de aumentos de densidad. Para incrementar la precisión de la técnica se utilizó nuevamente la técnica de codificación SDV.

Al detectarse una alta derivada positiva en la densidad, el sistema determina un “intervalo temporal” en el que se estima que ocurre la transición (Fig. 3.10.b, en gris). Si simultáneamente se detecta la caída en la DALPHA dentro del “intervalo temporal” estimado por la señal de densidad, el sistema da por reconocida una transición, estableciendo como instante de la misma el identificado en la señal DALPHA (57.08 segundos en la figura ejemplo). Esta simultánea coincidencia ocurrió en el 76% de los casos. El error medio (entre el instante determinado por los expertos y el sistema desarrollado) fue de 34 ms, con un error máximo de 89ms y una desviación estándar de 19.2 ms. Las estadísticas corresponden a 50 descargas en las cuales los expertos habían determinado el instante exacto de la transición.

---

<sup>18</sup> Pastillas congeladas de combustible que se lanzan al plasma a gran velocidad.

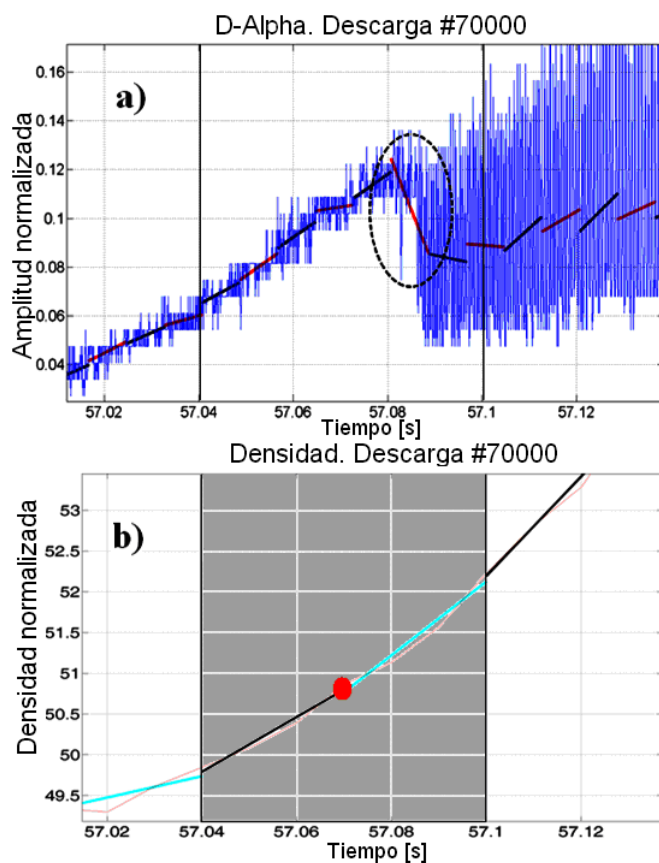


Fig. 3.10 Identificación combinada para el reconocimiento de transiciones.

- a) Identificación de la transición en una señal DALPHA. La primitiva de mayor longitud determina el instante de la transición. b) En una señal de densidad se define una zona (en gris) dentro de la que se estima que se encuentra el instante de la transición.

La implementación de estos métodos debe entenderse como una prueba del concepto y puede ser de gran utilidad cuando la base de datos no contiene una cantidad de ejemplos suficientes como para entrenar adecuadamente un sistema de aprendizaje. Con ellos, tal como se demuestra en otros trabajos [44, 45], es posible obtener mejores resultados.

### 3.4 CONCLUSIONES.

La aplicación de métodos de minería de datos para búsquedas automáticas en bases de datos masivas es indudablemente necesaria en fusión nuclear. En este capítulo se explica la aplicación de técnicas de reconocimiento morfológico de patrones dentro de señales. Estos métodos ya están disponibles desde los ordenadores del laboratorio internacional JET. Las técnicas son útiles para búsquedas de carácter general. Es decir,

para cuando un usuario está interesado en encontrar, dentro de una señal de referencia, alguna porción de interés y obtener de la base de datos estructuras semejantes a la seleccionada.

Para reconocer fenómenos físicos específicos se aprovechó e incluyó en las condiciones de búsqueda el conocimiento concreto sobre el comportamiento de las señales. Estos métodos *ad hoc* fueron aplicados, como pruebas de concepto, para analizar dos fenomenologías físicas específicas en plasmas de JET. En primer lugar se adaptó la técnica para realizar búsquedas capaces de detectar cortes en canales de temperatura. En segundo lugar se determinó el instante de la transición L-H mediante una múltiple búsqueda de patrones, combinando las formas morfológicas de las señales DALPHA y densidad integrada de línea de JET. Estas aplicaciones, al proveer resultados con gran velocidad (los tiempos de búsqueda son aproximadamente de 1ms por cada señal almacenada en la base de datos), son herramientas particularmente apropiadas para fusión, donde las bases de datos guardan cantidades cada vez mayores de información. Estos métodos proporcionan una solución rápida e intuitiva, proveen altas tasas de aciertos y además pueden ser extendidos a cualquier otra fenomenología cuyo comportamiento sea posible describir morfológicamente.





---

## **4 Extracción de modelos físicos mediante técnicas de aprendizaje: La transición L-H**

---

### 4.1.1 INTRODUCCIÓN.

En la Sección 3.3.2.2 del capítulo anterior se aplicaron técnicas de reconocimiento morfológico para la identificación automática de transiciones de modo L a modo H. Este tipo de herramientas de búsqueda puede emplearse para la rápida detección de fenomenología específica. Sin embargo, además del necesario desarrollo de herramientas que faciliten los procesos de búsqueda, la interpretación de la información surge como un tema fundamental. Una vez identificado cierto fenómeno físico es de gran utilidad construir automáticamente un modelo matemático que lo describa con sencillez y precisión. Estos modelos, generados a partir de los datos mediante sistemas de aprendizaje automático, tienen una gran importancia científica.

En el caso puntual de las transiciones L-H, para entender en profundidad la adecuada relación de parámetros que la producen, es necesario desarrollar modelos más complejos. Al tratarse de un fenómeno físico que ocurre no sólo en un Tokamak sino en todas las máquinas con esta configuración, es también deseable hallar una fórmula válida y general. Es decir, que el modelo sea escalable y aplicable a varias máquinas de características similares.

Al respecto, las leyes de escala, en la mayoría de las ramas de las ciencias, suelen ser expresadas en términos de productos de monomios con diferentes exponentes:

$$A = \frac{B^{kb} C^{kc} \dots D^{kd}}{H^{kh} I^{ki} \dots J^{kj}} \quad (4.1)$$

Donde las  $k$  representan valores constantes y las letras en mayúsculas magnitudes físicas.

En las próximas secciones se explicará cómo, mediante el uso de SVM, es posible desarrollar automáticamente modelos de la forma descrita anteriormente. Estas ecuaciones, extraídas mediante la aplicación de un sistema de aprendizaje, utilizarán las mismas magnitudes físicas que los modelos teóricos más aceptados para la determinación de las transiciones. Tanto los modelos extraídos mediante el uso de SVM como los teóricos se probarán con los datos provistos por una base de datos internacional, que incluye mediciones de los Tokamaks más importantes del mundo.

Las tasas de aciertos obtenidas con los modelos SVM en este estudio son superiores a las de los modelos teóricos. Esto no demuestra que los modelos deducidos mediante técnicas de aprendizaje puedan sustituirlos. Es la prueba de que dada una base de datos, son capaces de relacionar eficientemente los parámetros que describen comportamientos físicos y de que pueden ser utilizados como una herramienta de primera aproximación para el desarrollo o la confirmación de teorías.

#### 4.1.2 TRANSICIONES L-H.

Debido a los altos requerimientos pretendidos en el futuro dispositivo ITER, será necesario que este trabaje en modo H. Por lo tanto, una prioritaria cantidad de tiempo y esfuerzo en diferentes campañas experimentales de otros dispositivos se ha destinado a identificar el umbral de potencia necesario para superar el modo de bajo confinamiento. La serie de diferentes modelos teóricos existentes con el fin de estimar ese umbral de potencia consideran varias mediciones de parámetros del plasma. Los modelos suelen basarse en que las transiciones ocurren cuando se supera una temperatura electrónica ( $T_e$ ) crítica, ya que ésta puede relacionarse con la potencia externa aplicada al plasma.

Los modelos más aceptados suelen quedar expresados en forma de productos de monomios de la forma:

$$T_{ec} \propto p_1^a p_2^b \dots p_n^h \quad (4.2)$$

Donde  $T_{ec}$  es la temperatura umbral a superar para alcanzar la transición y los  $p_i$  representan los diferentes parámetros del plasma. La proporcionalidad indica que la ecuación debe ser multiplicada por alguna constante. Si la temperatura es mayor que  $T_{ec}$  el modelo estima que se ha alcanzado un modo de alto confinamiento.

#### 4.1.3 EXTRACCIÓN DE CARACTERÍSTICAS.

El reconocimiento de régimen en el caso de las transiciones consiste, básicamente, en un problema de clasificación bi-clase: dados unos parámetros, se pretende distinguir si el plasma se encuentra en modo L o en modo H. Si para el entrenamiento de un sistema basado en SVM se utiliza un kernel lineal, la frontera que separará los objetos según su clase quedará definida mediante un hiperplano de separación de forma:

$$k_a.A + k_b.B + k_c.C + \dots + k_d.D + cte = 0 \quad (4.3)$$

Aquí las  $k$  representan las constantes resultantes de la optimización calculada por SVM y las letras en mayúsculas las diferentes magnitudes físicas medidas. Los términos de esta expresión no poseen la forma deseada de la ec. (4.4). Sin embargo, si las magnitudes físicas son expresadas en forma logarítmica, es posible transformarlas en monomios de la forma de la ec. 4.1. dado que:

$$k_a \log(A) + k_b \log(B) + k_c \log(C) + \dots + k_d \log(D) + cte = 0 \quad (4.5)$$

$$\log(A^{k_a} . B^{k_b} . C^{k_c} \dots D^{k_d}) = -cte \quad (4.6)$$

$$A = k . B^{k_b'} . C^{k_c'} \dots D^{k_d'} \quad (4.7)$$

(nótese su semejanza con la ec. 4.2)

Donde:

$$k = -cte / k_a$$

$$k_b' = k_b / k_a$$

$$k_c' = k_c / k_a$$

$$k_d' = k_d / k_a$$

#### 4.1.4 BASE DE DATOS Y MODELOS TEÓRICOS.

Para comparar los resultados obtenidos mediante las ecuaciones derivadas de los clasificadores SVM con los de los modelos teóricos se utilizó la base de datos internacional [46]. El software público utilizado para la implementación del sistema SVM fue Spider [47]. Esta base de datos fue concebida explícitamente para estudios sobre el modo H e incluye señales validadas de la mayoría de los Tokamaks del mundo. De ella se extrajeron mediciones de 400 instantes en las que las operaciones se encontraban en el modo L y 400 en el modo H. El 70% de estos datos fueron empleados como entradas de entrenamiento para el clasificador SVM, reservando el 30% restante para la etapa de validación.

Los modelos teóricos considerados [48 - 53] son ampliamente aceptados por la comunidad científica internacional. Los propuestos por Chankin y col. [48, 49] asumen que los mecanismos de transporte del modo L están basados en turbulencias de ondas de deriva y que los efectos resistivos peliculares se consideran el mayor factor estabilizador. Kernel propone dos modelos detallados en [50]. Ambos se basan en inestabilidades en las ondas de deriva. Uno de ellos asume un transporte convectivo y el otro uno conductivo en la capa externa de la configuración magnética. En el modelo formulado por Rogister [51] también se consideran los efectos de las partículas neutras en el centro del plasma. La turbulencia en las ondas de deriva en una geometría toroidal es la base del modelo computacional propuesto en [52]. En [53], el mecanismo de supresión de turbulencia que dispara el modo H se supone vinculado al flujo cizalla generado por las pérdidas en las órbitas de los iones.

#### 4.1.5 RESULTADOS.

Los modelos derivados matemáticamente de las ecuaciones del hiperplano calculado con SVM se obtienen mediante el proceso de extracción de características detallado en la Sección 4.1.3. Los resultados se detallan en la Tabla 4.1. Allí puede advertirse que la cantidad de parámetros del plasma utilizados por los modelos generados SVM son los mismos que los modelos teóricos. Por fines comparativos, entonces, para el primer modelo de la tabla (Chankin) se entrenó el sistema SVM con las magnitudes medidas de campo toroidal ( $B$ ), radio mayor de la máquina ( $R$ ) y factor de seguridad ( $q$ ), las mismas variables empleadas en el modelo teórico de Chankin.

Modelo teórico	Ecuación	%	Ecuación del modelo SVM	%
Chankin	$T_{ec} \propto \left(\frac{\tau}{1+\tau}\right)^{3/16} \frac{(Rq)^{1/8} B^{1/2}}{A^{1/16}}$	87.9%	$T_{ec} = 311 * \frac{B^{1.17} R^{0.86}}{q^{0.28}}$	93.1%
Kernel et al. (Collisionless)	$T_{ec} \propto \frac{B^2}{(Rq)^{4/5} A^{3/5} n^{7/5}}$	91.3%	$T_{ec} = 821 * \frac{B^{1.14} R^{0.89}}{q^{0.36} n^{0.02}}$	93.4%
Kernel et al. (Collisional)	$T_{ec} \propto \frac{(Rq)^{2/15} n^{1/15} B^{2/5}}{A^{1/15}}$	83.6%		
Rogister et al.	$T_{ec} \propto \frac{AqRn^2}{aB^2}$	50%	$T_{ec} = 4026 * \frac{B^{1.12} a^{0.90}}{q^{0.4} R^{0.61} n^{0.02}}$	94.7%
Scott et al.	$T_{ec} \propto \frac{\tau B^2}{nA(1+\tau)}$	90.5 %	$T_{ec} = 1090 * \frac{B^{1.67}}{n^{0.02}}$	91%
Shaing & Crume	$T_{ec} \propto \frac{\tau R^{5/4} (qn)^{1/2}}{a^{3/4}}$	63.3%	$T_{ec} = 11854 * \frac{a^{1.28}}{q^{0.19} R^{0.9} n^{0.17}}$	89.3%

Tabla 4.1 Los resultados de los diferentes modelos teóricos son comparados con los de las ecuaciones extraídas automáticamente del clasificador SVM desarrollado.

Estos últimos, en todos los casos, proveen una mayor tasa de aciertos.

Consecuentemente, se seleccionó el mismo conjunto de parámetros en los sistemas SVM que los utilizados en los otros modelos teóricos a comparar.

Una vez obtenidas las ecuaciones finales en forma de monomios, las tasas de acierto se calculan de la misma forma para los modelos teóricos y los de SVM:

- Se introducen en las ecuaciones los valores de los parámetros del plasma.

- La ecuación determinará si se ha superado la temperatura crítica y por lo tanto si se está en modo L o H.
  - Si la predicción es correcta (en las bases de datos se incluye si el régimen es L o H) se la considera un acierto. Si es incorrecta, como un error.
- Se repite el mismo procedimiento para todos los valores de la base de datos y se calculan las tasas totales de acierto.

Para la correcta interpretación de los resultados son necesarias, además, algunas observaciones:

- Al contrario de las ecuaciones derivadas con SVM, los modelos teóricos proveen sólo una proporcionalidad entre los valores críticos de temperatura electrónica y las demás cantidades de escala.
- En consecuencia, para cada uno de los modelos teóricos fue necesario calcular empíricamente un factor de proporcionalidad. Para ello se programó un algoritmo que iterativamente asignaba un valor de proporcionalidad. Después de cada iteración, se calculaba la tasa de aciertos. El valor de la constante que maximizaba la tasa de aciertos fue el utilizado en cada caso.
- Esto se debe a que tales ecuaciones requieren que este parámetro sea calculado para cada base de datos en las que se empleen. En cualquier caso, conforme a la Tabla 4.1, incluso con este factor calculado específicamente para cada caso, el rendimiento de los modelos teóricos es sustancialmente más bajo que el de los desarrollados mediante SVM (que oscilan entre el 89% y el 95% de aciertos).

Cualquier modelo SVM se construye a través de los datos con los que el sistema fue entrenado. La base de datos internacional de transiciones L-H utilizada incluye información de una gran cantidad de dispositivos. Sin embargo, no es posible garantizar que los modelos derivados de SVM obtengan tasas de aciertos similares si se prueban en máquinas diferentes a las incluidas en la base de datos internacional. Las ecuaciones extraídas de SVM pueden proporcionar una primera aproximación sobre las relaciones

entre los parámetros almacenados pero no garantizan una descripción absoluta de la física de las transiciones.

Un artículo en el que se describe en detalle este trabajo y en el que además se generan modelos mediante otra técnica de aprendizaje supervisado (redes neuronales) ha sido enviado recientemente para su publicación [54].

## **4.2 CONCLUSIONES.**

En este capítulo se implementó SVM con el objetivo de generar modelos capaces de describir las transiciones entre regímenes de operación. Para ello se utilizaron las propiedades de SVM para clasificación lineal en problemas bi-clase. La solución hallada en estos casos consiste en la ecuación de un hiperplano que depende de las variables de entrenamiento.

Para adecuar las soluciones a la forma de monomios se introdujeron a SVM los logaritmos de las magnitudes empleadas. De esta forma fue posible adecuar las soluciones obtenidas con SVM para que fueran comparables con las de los modelos teóricos. SVM calcula la mejor relación entre magnitudes físicas para describir este fenómeno físico de gran interés en la comunidad científica relacionada a la fusión nuclear.

Los resultados, basados en las tasas de acierto obtenidas, demuestran su certeza en la predicción del régimen en la que el plasma se encuentra. Estos porcentajes de acierto son superiores, en todos los casos, a los obtenidos mediante los modelos físicos más comúnmente aceptados.





---

## **5 Análisis diferido de disrupciones en JET**

---

### **5.1 INTRODUCCIÓN Y ESTADO DEL ARTE.**

La configuración Tokamak se presenta como una seria opción para el desarrollo de un futuro reactor. Sin embargo se ve sometida constante e inevitablemente a pérdidas abruptas del confinamiento del plasma llamadas disrupciones [55]. Estas repentinas inestabilidades que causan el fin de la descarga, además de afectar a la continuidad de la operación del dispositivo, pueden constituir un serio factor de riesgo para su integridad. Durante una disrupción, en períodos de tiempo del orden de milisegundos, el plasma excede sus límites operacionales con la resultante pérdida de energía y corriente. Durante la primera fase de la disrupción es posible que se transfieran altísimas cargas térmicas a los componentes de la primera pared de la máquina. Posteriormente, grandes corrientes son inducidas en la cámara de vacío y en las estructuras aledañas provocando fuerzas capaces de causar daños considerables. Actualmente se ha probado que su aparición durante la operación es inevitable, especialmente en configuraciones de alto rendimiento. La inminente necesidad de eludir el fenómeno mediante su temprana detección será incluso mayor si se tiene en cuenta el próximo Tokamak ITER. En ITER se planean obtener plasmas con mayores corrientes y de tales energías que una sola disrupción podría provocar daños severos en varias partes del dispositivo, desde la erosión de los componentes expuestos al plasma (aquellos que recubren internamente la

cámara de vacío y resto de sistemas próximos a los flujos de partículas y energía procedentes del plasma) a daños en su estructura.

La caracterización física de las disrupciones para su posible predicción y control es extremadamente compleja. La gran cantidad de variables involucradas en el fenómeno y la relación no lineal entre ellas han hecho imposible hasta el momento desarrollar un modelo físico lo suficientemente fiable capaz de detectar satisfactoriamente este repentino comportamiento. Por ello en los últimos 15 años varios sistemas de aprendizaje (fundamentalmente basados en distintas Redes Neuronales Artificiales o Máquinas de Vectores Soporte) han sido utilizados como una aproximación alternativa a la detección del fenómeno. Los sistemas de aprendizaje son aplicables dado que las disrupciones pueden considerarse un problema de clasificación. Las clases a identificar en una descarga serán si esta disrumpirá o no.

Muy pocos de los trabajos desarrollados hasta el momento estudian la extracción de características de las señales empleadas para el entrenamiento de los sistemas de aprendizaje. El apropiado procesado de los parámetros del plasma ayuda a evidenciar con mayor claridad el fenómeno y por lo tanto debe entenderse como una herramienta necesaria para la obtención de altas tasas de aciertos en cualquier sistema de clasificación. Entre los estudios que sí abordan tanto la adecuada extracción de características como la identificación de comportamientos disruptivos se puede destacar el desarrollado por Cannas y col. [56], publicado en el año 2004. Mediante el entrenamiento de un sistema de redes neuronales, se demuestra que el clasificador es capaz de diferenciar entre descargas disruptivas y no disruptivas. Como método de extracción de características y de reducción de dimensionalidad aplica mapas autoorganizados de Kohonen a porciones de nueve señales para el entrenamiento del sistema clasificador. Estas señales son: la corriente del plasma, mode lock, energía total irradiada, densidad electrónica, potencia total de entrada, inductancia interna del plasma, derivada temporal de la energía diamagnética almacenada, factor de seguridad al 95% y beta poloidal. El sistema logra una efectividad del 95% en la distinción entre partes de descargas disruptivas y no disruptivas. Posteriormente, en el año 2006, en el Tokamak alemán ASDEX Upgrade, se realizó un estudio semejante al anterior en cuanto a la extracción de características pero empleando en esta ocasión siete parámetros del plasma [57]. Los parámetros son: mode lock, densidad electrónica,

potencia total de entrada, inductancia interna del plasma, factor de seguridad al 95%, beta poloidal y potencia de red (energía total irradiada/ potencia total de entrada). Con esas señales se efectúa el análisis únicamente sobre descargas disruptivas, obteniéndose un 90% de correctas identificaciones. Cabe destacar que en [57] las descargas eran analizadas de principio a fin (si bien sólo disruptivas, lo que no permite su aplicación en tiempo real). Una técnica desarrollada con lógica difusa en JET (año 2008) por Murari y col. [58] debe ser también mencionada. Este clasificador basado en reglas no alcanza las tasas de aciertos de los predecesores. Sin embargo, presenta un innovador análisis, mediante el CART, de la importancia relativa de las diferentes señales medidas para la identificación del fenómeno.

Algunas disrupciones ocurren considerablemente más rápido que otras y por lo tanto, como se hablará más adelante, pueden clasificarse según su tipo. En los tres trabajos anteriores, y en la amplia mayoría de los publicados hasta el momento, los tipos de disrupciones de difícil detección<sup>19</sup> son descartados tanto en la etapa de desarrollo de los clasificadores como en las de pruebas y consecuentemente de las estadísticas finales calculadas. No son considerados los experimentos que incluyen corrientes menores a 1,5 MA, arguyéndose que en tales condiciones la severidad de los efectos de las disrupciones decae considerablemente.

Es especialmente importante enfatizar que ninguno de los trabajos mencionados es aplicable a condiciones de tiempo real debido a que no analizan la completa evolución temporal de cada descarga (disruptivas y no disruptivas), de principio a fin. El análisis, en cambio, se acota a periodos (o pulsos) específicos. Dichos estudios deben considerarse como un paso necesario para el desarrollo de sistemas aplicables en tiempo real, pero no como una solución definitiva.

En este capítulo se selecciona un conjunto válido de señales para la adecuada identificación de disrupciones. De ellas se extraen características apropiadas para el entrenamiento de un sistema de aprendizaje que permita clasificar las descargas en disruptivas o no disruptivas [59]. A diferencia de la gran mayoría de trabajos publicados previamente, en el presente no se descarta ningún tipo de disrupción ni se obvian

---

<sup>19</sup> La evolución de los precursores que llevan a una descarga a disrumpir varía en sus características y tiempos. En algunos casos la detección del fenómeno es considerablemente más difícil que en otros.

descargas que cumplan condiciones particulares, como no ser estacionarias o tener una corriente toroidal menor que 1,5 MA. Lo que se pretende es dar el primer paso para, como se explica y desarrolla en el siguiente capítulo, construir un modelo capaz de predecir con efectividad todo tipo de disrupciones durante la operación de la máquina en tiempo real.

Este capítulo está esquematizado de la siguiente manera. Inicialmente se introducen las fases y los tipos de disrupciones más comunes (Sección 5.1.1). La Sección 5.2 profundiza en el estudio sobre la extracción de características apropiada para la detección del fenómeno. El proceso de extracción de características se verifica en la Sección 5.2.3.2 mediante la implementación de técnicas para la visualización de datos de alta dimensionalidad. En la Sección 5.3 se detalla el procedimiento de entrenamiento de un sistema basado en SVM capaz de distinguir entre porciones de descargas disruptivas y no disruptivas, demostrando que a medida que un experimento se aproxima al instante de la disrupción el fenómeno se hace más evidente.

Finalmente, en la Sección 5.4 se resumen y comentan los resultados obtenidos.

### 5.1.1 FASES Y TIPOS DE DISRUPCIONES.

#### 5.1.1.1 Fases.

Las disrupciones pueden ocurrir repentinamente, prácticamente sin comportamientos anormales obvios o después de periodos del orden de cientos de milisegundos, en los cuales diferentes inestabilidades se potencian hasta terminar abruptamente la descarga. De todas maneras está generalmente aceptado que el fenómeno evoluciona en una serie de fases [55, 60].

A la primera fase se la llama “*pre-precursor*” o “*evento inicial*”. Allí la energía radiada o la densidad se incrementan debido a anomalías externas (como por ejemplo fallos mecánicos o de la operación) o internas (comportamientos imprevistos) del plasma, aunque esas variaciones suelen ser sutiles y por lo tanto la posibilidad de que sean detectadas es mínima.

Durante la “*fase de precursores*”, las irregularidades iniciadas en la fase anterior alcanzan un punto crítico determinado por la aparición de crecientes y detectables inestabilidades MHD (MagnetoHidroDinámicas). La importancia de esta fase en la

detección temprana es fundamental: el fenómeno comienza a ser evidente y además, si es detectado con la suficiente rapidez, es posible iniciar acciones de control para evitarlo o al menos atenuar su efecto. Teóricamente estas inestabilidades comienzan como islas magnéticas que rotan alrededor del eje magnético del plasma y que crecen de tamaño exponencialmente en pocas decenas de milisegundos hasta llevar al plasma a la disrupción.

A la fase siguiente se la conoce como “*fase rápida*”. En esta, en tiempos del orden de milisegundos, el perfil radial de corriente se aplanan y la temperatura en el centro del plasma cae con rapidez.

A la fase final se la llama la “*fase de apagado*”. Durante el “*apagado térmico*” la temperatura electrónica cae bruscamente dejando al plasma en un estado altamente resistivo. El calentamiento óhmico resultante llega a ser del orden del Gigawatt. El aumento de la resistividad deriva parte de la corriente toroidal, durante la fase de “*apagado de corriente*”, hacia la cámara de vacío. Estas corrientes inducen momentos de gran magnitud en las estructuras colindantes a la primera pared, las cuales sumadas a las cargas térmicas transferidas en escalas de tiempo menores a 1 ms pueden causar daños severos al dispositivo.

#### 5.1.1.2 Tipos de disrupciones.

Las causas de las disrupciones que se producen como consecuencia de comportamientos anómalos del plasma no son totalmente comprendidas todavía. Sin embargo, existen diferencias en su desarrollo inicial y en las fases finales (en las que se pierde la energía) en virtud de las cuales puede establecerse una clasificación dependiendo de su evolución. Los tipos más comúnmente aceptados son los mencionados a continuación:

Las disrupciones de “*límite de densidad*” [61] pueden ser descritas como una inestabilidad que comienza cuando la densidad se eleva más allá de un umbral o cuando existen impurezas en el plasma. La energía radiada crece y consecuentemente se puede medir una caída de la temperatura. Además, el perfil de corriente se encoje como consecuencia del aumento de la resistencia del borde del plasma. Por propiedades

inductivas del plasma, el *apagado de corriente* produce corrientes de halo<sup>20</sup> transfiriendo grandes fuerzas a la cámara de vacío. Este tipo de disrupciones es el más común en JET.

Otro tipo muy frecuente es el “*mode lock*” [61]. Debido a gradientes de temperatura o campos eléctricos externos se desarrollan inestabilidades MHD produciendo islas magnéticas, las cuales crecen hasta terminar en una disrupción.

Las impurezas en el plasma son las causantes de las disrupciones llamadas “*high radiated power*” [58]. La evolución de tales disrupciones se asemeja a las de límite de densidad pero con la diferencia que el incremento de energía radiada se produce ligeramente antes que el aumento en la densidad.

Las disrupciones de “*límite vertical*” o “*desplazamiento vertical*” [61] ocurren normalmente en configuraciones de plasmas elongados. Muy rara vez el sistema de control permite este tipo de condiciones y por lo tanto en general se producen cuando los campos de estabilización vertical son desactivados expresamente para experimentos puntuales.

También transiciones de regímenes de alto a bajo confinamiento (“*H-mode/L-mode disruptions*”) [58] a altas densidades pueden desarrollar inestabilidades causantes de disrupciones.

Las disrupciones de “*Alta Beta*” [61] están vinculadas a incrementos de la presión del plasma que desencadenan una inestabilidad MHD. De hecho,  $\beta$  es un parámetro MagnetoHidroDinámico (MHD) que mide la presión de confinamiento y por lo tanto se la considera directamente relacionada con los límites de estabilidad. En JET estos límites no se alcanzan en condiciones de operación normal y por lo tanto este tipo de disrupciones no es frecuente.

Se ha demostrado que para superficies magnéticas, en general con valores bajos del factor de seguridad  $q$ <sup>21</sup> y en particular aquellas con  $q=2$  juegan un papel

---

<sup>20</sup> Aquellas que circulan en el plasma frío fuera de la última superficie magnética cerrada.

<sup>21</sup> Número de vueltas que las líneas de campo magnético helicoidales dan en dirección toroidal por cada vuelta en dirección poloidal. Para evitar inestabilidades es necesario que el factor sea superior a 1.

fundamental en la estabilidad del plasma. La superficie  $q=2$  suele volverse inestable y por lo tanto causar las disrupciones por bajo factor de seguridad [62] (“*Low q*”).

Finalmente, las ITB (“*Internal Transport Barriers*”) [58], son modos que suelen aparecer en las configuraciones etiquetadas como “*escenarios avanzados*” y su detección es particularmente complicada dado el breve margen entre la fase de *precursor* y la de *apagado de corriente*.

## 5.2 EXTRACCIÓN DE CARACTERÍSTICAS Y PROCESADO DE LAS DESCARGAS.

### 5.2.1 BASE DE DATOS Y PRE-PROCESADO.

Para realizar el estudio de extracción de características se dispuso de 440 descargas de JET de las cuales 220 terminaban en una disrupción. Para asegurar la fiabilidad de cada una de las señales involucradas en cada experimento, un grupo de expertos determinó previamente el instante exacto de la disrupción en el caso de las descargas disruptivas.

Las señales son adquiridas a diferentes frecuencias, dependiendo del diagnóstico utilizado para cada medición. Para su tratamiento computacional y su debida sincronización, el conjunto seleccionado de parámetros del plasma debe remuestrearse con el fin de que cada señal tenga el mismo número de muestras. Cabe añadir que debido a las diferencias de varios órdenes de magnitud en las amplitudes adquiridas para las diferentes señales es necesario primeramente aplicar una función de normalización. Esta garantiza que cada uno de los parámetros involucrados en el estudio esté limitado entre los valores mínimos y máximos de 0 y 1 respectivamente, sin perder las magnitudes relativas. Ésta se puede expresar como:

$$\text{Señal Normalizada} = \frac{\text{Señal} - \text{Min}}{\text{Max} - \text{Min}} \quad (5.1)$$

donde Min y Max representan respectivamente los valores máximos y mínimos de cada señal para la base de datos de las 440 descargas.

En consecuencia, las diferentes señales que se emplean quedan representadas por amplitudes comparables, evitando que los sistemas utilizados para la clasificación

asignen un mayor peso debido a la magnitud de la señal. Por lo tanto, la normalización efectuada proporciona un mismo peso relativo a todas las señales.

### 5.2.2 SELECCIÓN DE LAS SEÑALES MÁS RELEVANTES Y USO DE VENTANAS TEMPORALES.

Para el caso de las interrupciones, sólo un número reducido de la gran cantidad de señales generadas contiene información útil para el estudio del fenómeno. Es imprescindible seleccionar un conjunto compacto que aporte toda la información relevante y la menor cantidad de datos redundantes. Este problema fue abordado en trabajos previos [56, 59]. En la mayoría de ellos en la elección de unas 9 señales, detalladas en la parte no sombreada de la Tabla 5.1. Estas señales serán las de referencia de aquí en adelante, con ciertas excepciones: tres derivadas temporales de algunas de ellas (resaltadas en gris) debido a que se ha demostrado que pueden ayudar a evidenciar comportamientos disruptivos [58, 63] y la señal de posición vertical del plasma. Ésta última (resaltada en rojo en la Tabla 5.1) no suele estar incluida en otros estudios debido a que fundamentalmente es útil para reconocer interrupciones rápidas, como las de límite vertical. Este tipo de interrupciones, por ser generalmente provocadas intencionalmente para estudios específicos y por su alta dificultad de detección, no suelen ser incluidas ni para el entrenamiento de sistemas predictores ni para el posterior test de los mismos. Sin embargo, el fin de este trabajo es el de crear un sistema de detección general que no excluya ningún tipo de interrupción, tanto en la etapa de entrenamiento como en la de pruebas, y por lo tanto este parámetro también fue considerado.

La importancia de cada una de las señales para evidenciar comportamientos disruptivos varía a medida que la descarga se aproxima al fenómeno, tal como se demuestra en varias publicaciones [69, 70, 71]. Esto se demuestra al introducir las 9 señales a un sistema de clasificación CART. La relevancia de las mismas para la detección de interrupciones cambia de acuerdo con lo expresado en la Tabla 5.2. Allí, en la primera fila se representan los intervalos de tiempo (en milisegundos) analizados. En estos casos el tiempo 0 corresponde al instante de la interrupción.



Nombre de la señal	Unidades
1. Corriente del plasma.	$A$
2. Beta poloidal.	
3. Derivada temporal de 2.	$s^{-1}$
4. Amplitud de “mode lock”.	$T$
5. Factor de seguridad a radio 0.95.	
6. Derivada temporal de 5.	$s^{-1}$
7. Potencia total introducida.	$W$
8. Inductancia interna del plasma.	
9. Derivada temporal de 8.	$s^{-1}$
10. Posición vertical del plasma.	$m$
11. Densidad.	$m^{-3}$
12. Derivada temporal de la energía diamagnética acumulada.	$W$
13. Potencia de red (7 menos energía total radiada).	$W$

Tabla 5.1 Listado de las señales seleccionadas.

Para poder realizar este análisis (y también otros posteriores, como se verá en las siguientes secciones) es necesario utilizar franjas temporales específicas de las descargas que llamaremos “ventanas temporales”, las cuales representan un intervalo de tiempo. Cada uno de estos intervalos (cuya obtención se detalla en las secciones posteriores) condensará la información relevante al fenómeno disruptivo de un período de descarga. De ahora en más nos referiremos a ellos como ventanas temporales o vectores de características.

Todo análisis de los comportamientos de las descargas se realizará sobre estas ventanas temporales correspondientes a periodos vinculados al instante de la disrupción, como por ejemplo el comprendido en el intervalo [-60, -30] milisegundos antes de la disrupción.

Para fines comparativos (ya que el intervalo [-60, -30] queda definido perfectamente para descargas disruptivas) se deriva la necesidad de estimar un instante “equivalente” al de la disrupción para las descargas **no disruptivas**. Ese tiempo se definió como el de los 7 segundos posteriores a la formación del punto X, porque estadísticamente resulta ser el instante más probable para que ocurra una disrupción en descargas de JET [58]. Existe en JET una señal específica que indica el tiempo en el que se forma el punto X. De este modo, las ventanas correspondientes al intervalo [-60, -30]

[-320, -280] ms		[-260, -220] ms		[-200, -160] ms		[-140, -100] ms		[-80, -40] ms	
Variable	Relev.	Variable	Relev.	Variable	Relev.	Variable	Relev.	Variable	Relev.
$P_{net}$	100	$dW_{dia}/dt$	100	$dW_{dia}/dt$	100	$dW_{dia}/dt$	100	$dW_{dia}/dt$	100
$Ipla$	94,96	$P_{net}$	51,00	$Dens$	30,63	$\beta_p$	23,87	$Ipla$	6,65
$dW_{dia}/dt$	43,71	$\beta_p$	47,75	$\beta_p$	27,25	$l_i$	9,72	$Loca$	4,55
$Loca$	23,80	$l_i$	28,44	$l_i$	9,53	$Loca$	7,21	$l_i$	1,50
$\beta_p$	18,32	$Ipla$	19,36	$P_{net}$	9,19	$Dens$	4,50	$Dens$	1,34
$Dens$	16,39	$Dens$	13,36	$Ipla$	6,27	$q_{95}$	2,40	$\beta_p$	0,69
$q_{95}$	12,97	$Loca$	10,62	$Loca$	4,27	$Ipla$	1,24	$P_{net}$	0,44
$l_i$	6,49	$P_{inp}$	5,58	$q_{95}$	3,80	$P_{net}$	0,00	$P_{inp}$	0,0
$P_{inp}$	0,0	$q_{95}$	3,40	$P_{inp}$	0,0	$P_{inp}$	0,00	$q_{95}$	0,0

Tabla 5.2<sup>22</sup> Importancia relativa de cada señal (Relev.) asignada por el CART.

Este método fue utilizado para cada una de las 9 señales (en 4 diferentes ventanas temporales de 40 ms).

ms de las descargas disruptivas se compararán con las equivalentes [-60, -30] ms antes del “tiempo equivalente” en descargas no disruptivas.

La duración temporal de las ventanas es otro factor a ser determinado. En estudios previos estas oscilan habitualmente entre los 20 ms y los 40 ms. La elección final de fijar la duración de las ventanas en 30 ms en esta Tesis no está dada sólo por ser un valor razonablemente intermedio entre los más comúnmente empleados sino también por el problema a tratar. En JET, los actuadores disponibles necesitan un tiempo mínimo de 30 ms para efectuar acciones de mitigación una vez se detectan disrupciones y por lo tanto este tiempo resulta adecuado para la duración de las ventanas (la longitud temporal exacta de una ventana será el máximo margen para detectar y mitigar un comportamiento disruptivo).

<sup>22</sup> Extraída de [69].

### 5.2.3 EXTRACCIÓN DE CARACTERÍSTICAS.

#### 5.2.3.1 Introducción.

La extracción de características está relacionada con el procesado de las ventanas temporales para cada una de las 13 señales consideradas. El propósito es reducir la información redundante y resaltar la que evidencie comportamientos disruptivos.

Para abordar el problema se utilizaron diferentes fórmulas de procesado y visualización. Ellas posibilitan distinguir ventanas correspondientes a comportamientos disruptivos y no disruptivos. El procedimiento de extracción de características tuvo bases empíricas. Fue guiado por conocimientos generales sobre la evolución de las disrupciones y su consiguiente efecto en los diferentes parámetros del plasma, como por ejemplo variaciones repentinas de amplitud en algunas señales. La etapa de preprocesado (que consiste en el remuestreo y la normalización) fue empleada como base común para todas las descargas. En sí, los diferentes procesos de extracción de características (Fig. 5.1) consistieron en utilizar segmentos de 30 ms de las señales remuestreadas y normalizadas y aplicarles una amplia gama de técnicas como se detallará en la siguiente Sección.

#### 5.2.3.2 Vectores de características y sus representaciones visuales.

Inicialmente, y siguiendo las pautas empíricas mencionadas anteriormente, se consideraron 2 procesos diferentes de extracción de características. Ambos procesos se aplicaron a la misma ventana temporal: [-60, -30] milisegundos antes de la disrupción para las 220 descargas disruptivas y las 220 no disruptivas.

Llamemos conjunto de vectores de características a cada uno de los procesos de extracción que se detallarán a continuación. En este caso cada conjunto contendrá 440 vectores de características correspondientes a la ventana [-60, -30] milisegundos antes de la disrupción.

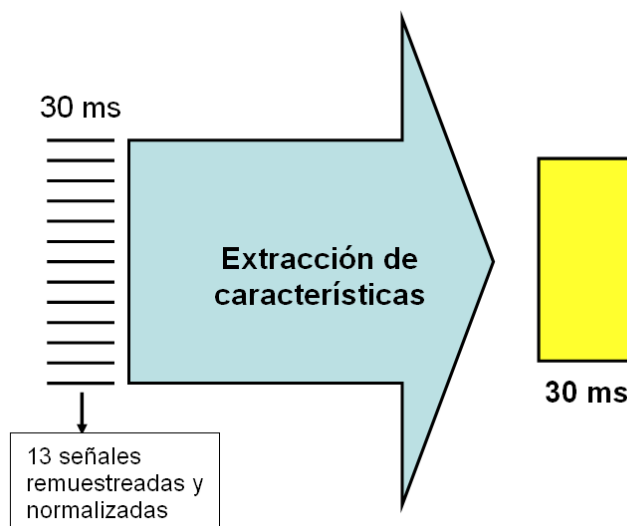


Fig. 5.1 Esquema del proceso de extracción de características.

Diferentes técnicas son aplicadas a ventanas de 30 ms correspondientes a las trece señales seleccionadas para crear un vector de características (representado por el rectángulo amarillo).

El primer conjunto de vectores se obtiene mediante la concatenación de las muestras preprocesadas (remuestreo y normalización) de las 13 señales. Contendrá entonces 440 vectores de características, cada uno de ellos con la información de 30 ms de una descarga.

Análogamente se obtiene el segundo conjunto de vectores de características. El proceso consiste, en este otro caso, en la aplicación de la transformada de Fourier a cada señal (reteniéndose el módulo del espectro y quitándole la componente continua). Finalmente, estos 13 espectros son concatenados.

Cada vector de características de cada conjunto puede considerarse como un objeto de alta dimensionalidad. Los del primer conjunto tendrán una dimensionalidad 390 (30 muestras por 13 señales). Los del segundo, 182 (14 muestras por 13 señales). Un buen proceso de extracción de características debería resaltar en mayor medida las diferencias entre los objetos pertenecientes a cada clase. Visualmente esto podría traducirse en la formación de 2 agrupaciones. Una mayor separación entre los grupos significará entonces una mejor extracción de características.

Para poder representar los conjuntos de vectores de alta dimensionalidad se recurrió a dos de los recursos detallados en el capítulo 2. El primero consistió en visualizar los vectores de alta dimensionalidad proyectándolos en un plano mediante el

pseudo Grand Tour (2 representaciones en la parte superior de la Fig. 5.2). El segundo, en reducir la dimensionalidad de los vectores de características mediante las técnicas de Mapas Topográficos Generativos (2 representaciones en la parte inferior de la Fig. 5.2).

En las representaciones, sobre todo en la del pseudo Grand Tour, se evidencian con mayor claridad los agrupamientos de las descargas disruptivas y no disruptivas para el segundo conjunto de vectores de características.

Estos recursos gráficos ayudan a comprender que en el dominio de las frecuencias la diferenciación entre ventanas disruptivas y no disruptivas es más evidente. Por lo tanto, los siguientes pasos para refinar el proceso de extracción de características se efectuaron en el dominio de las frecuencias. El objetivo consiste en continuar aplicando técnicas que:

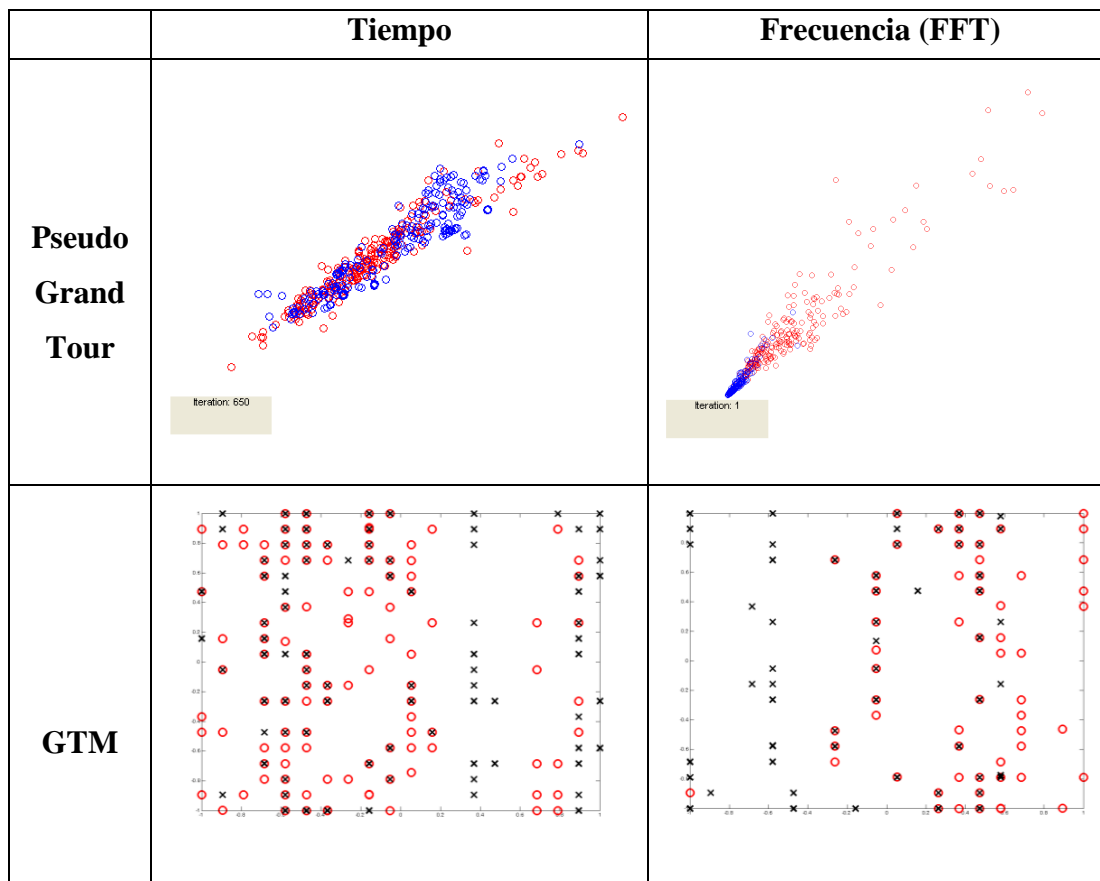


Fig. 5.2 Visualizaciones bidimensionales de dos conjuntos de vectores de características (mediante PsGT y MTG).

Se evidencian con mayor claridad las diferencias entre ventanas pertenecientes a descargas no disruptivas (circunferencias azules en pseudo Grand Tour y cruces negras en GTM) cuando se trabaja en el dominio de las frecuencias.

- 1- Reduzcan la dimensionalidad, perdiendo la mínima cantidad de información relevante relacionada con interrupciones.
- 2- Evidencien el fenómeno con mayor claridad.
- 3- Sean utilizables en un sistema de aprendizaje para que éste clasifique automáticamente entre descargas disruptivas y no disruptivas.

Para la visualización de las posteriores extracciones de características se eligió el PsGT, dado que sus proyecciones muestran con mayor claridad la separación entre los agrupamientos.

Todos los procedimientos de extracción de características probados fueron aplicados a las ventanas temporales de las señales siguiendo los siguientes pasos:

- 1- Preprocesado (remuestreo y normalización).
- 2- Aplicación de la FFT (calculándose el módulo y eliminándose la componente de continua) a cada uno de los segmentos de señales.
- 3- Aplicación de uno de los siguientes algoritmos al espectro obtenido:
  - PCA (2 ó 3).
  - SVD (2).
  - MDS (2).
  - Cálculo del valor medio (1).

Desviación estándar (1).

- 4- Los tres pasos anteriores se aplican a cada una de las 13 señales seleccionadas.
- 5- El quinto y último paso consiste en la concatenación de los valores obtenidos para cada señal, creando así el vector de características.

Los números entre paréntesis indican la cantidad de valores resultantes de cada operación aplicados a cada ventana temporal de una señal. En algunos casos, como PCA, se indica más de un valor dependiendo de la cantidad de componentes retenidas.

En la Fig. 5.3, se representan los diferentes vectores de características mediante el PsGT. En esta instancia la herramienta de visualización no es suficiente para demostrar de forma clara cuál de los vectores mejora la distinción de clases.

### 5.3 IMPLEMENTACIÓN DE MÉTODOS SUPERVISADOS.

#### 5.3.1 VALIDACIÓN CRUZADA.

El fin de la extracción de características es proveer las mejores entradas posibles a un sistema de clasificación. El objetivo es entrenar un sistema capaz de distinguir, para las diferentes ventanas temporales, entre comportamientos disruptivos y no disruptivos.

Una vez más el sistema de aprendizaje elegido fue SVM, especialmente adecuado para la clasificación y cuya facilidad de aplicación, buen rendimiento y disponibilidad en código libre ya fue mencionado en capítulos anteriores. El software público utilizado fue Spider [47]. Para cada uno de los casos de la Fig. 5.3. se probaron diferentes kernels obteniéndose, en general, los mejores resultados con los RBFs.

Para aumentar la fiabilidad de los resultados obtenidos (teniendo en cuenta que el número de descargas no es lo suficientemente amplio) se aplicó el algoritmo conocido como “n fold cross validation”, normalmente traducido como validación cruzada con n pliegues, con un  $n=4$ . Este método consiste en dividir el conjunto de datos aleatoriamente en los n grupos determinados. Un grupo se reserva para pruebas y mediante los restantes  $n-1$  se entrena un modelo.

El modelo se utiliza para la predicción de los resultados del grupo de pruebas. El proceso se repite n veces considerando en cada iteración un grupo de pruebas diferente. Se realizan n iteraciones y el porcentaje final de acierto se calcula como la media de esos valores.

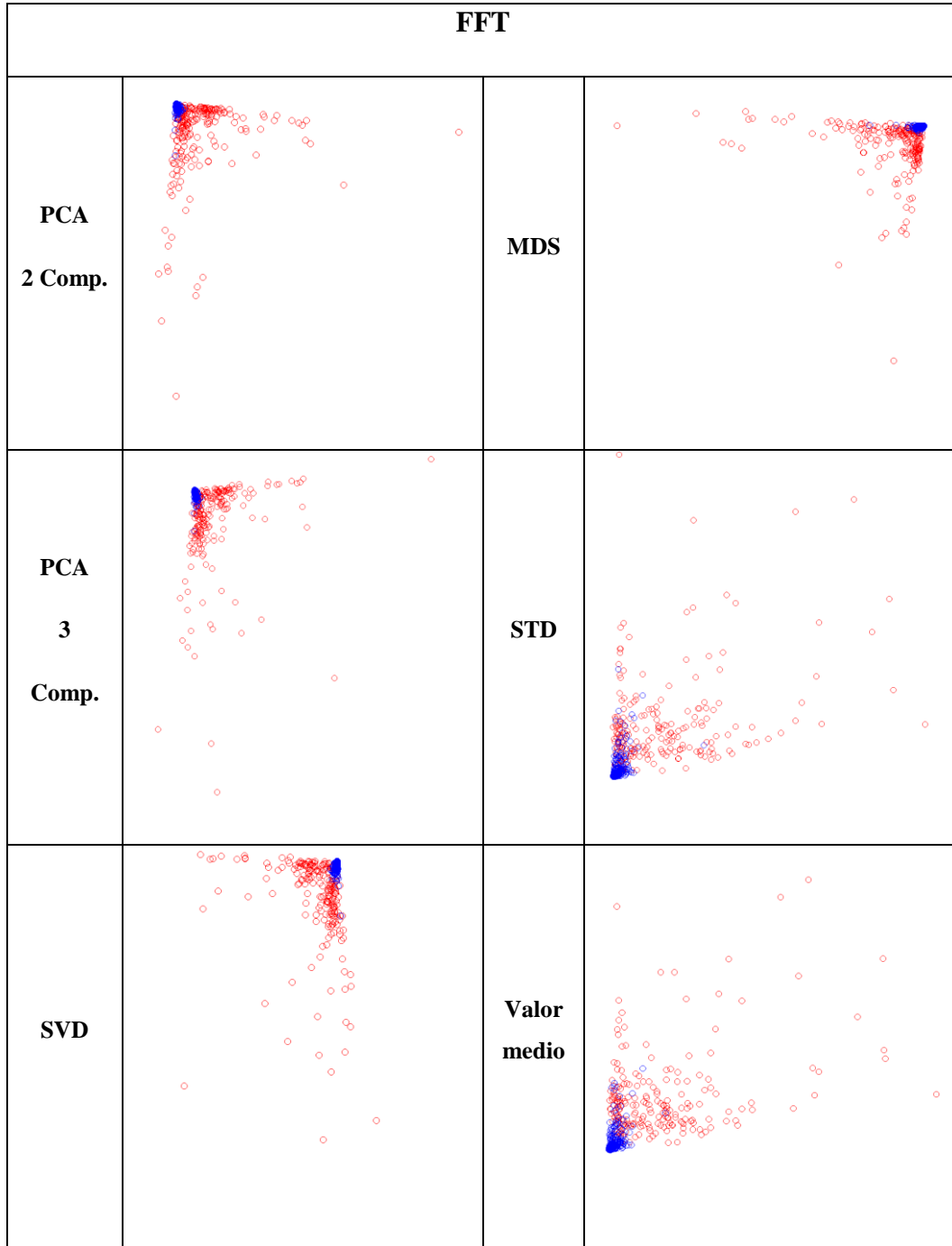


Fig. 5.3 Representación de seis vectores de características.

*En este punto las proyecciones no demuestran con claridad cuál vector mejora la diferenciación entre comportamientos disruptivos y no disruptivos.*

### 5.3.2 RESULTADOS.

Las tasas de acierto calculadas mediante validación cruzada, para diferentes intervalos de 30 ms, se muestran en la Fig. 5.4. Allí pueden observarse cuatro de las



trazas obtenidas: tres de ellas correspondientes a los vectores de características detallados en la Sección 5.2.3.2 y una con los datos normalizados y concatenados (primer vector de características representado en la Fig. 5.2).

Las trazas correspondientes a los vectores de características obtenidos mediante PCA (2 componentes), SVD y MDS no fueron pintadas. Estas se solapan en varios puntos de la curva de PCA con 3 componentes, sin aportar ninguna información de importancia y dificultando la comprensión de los resultados.

En la figura se destacan los porcentajes de acierto obtenidos mediante la utilización del siguiente conjunto de vectores de características: la desviación estándar del módulo de la transformada de Fourier de cada ventana (etiquetada como “desviación estándar”). Con esas características, las tasas son superiores al 94% hasta los ~180 ms antes de la interrupción y alcanzando un máximo del 98.3%.

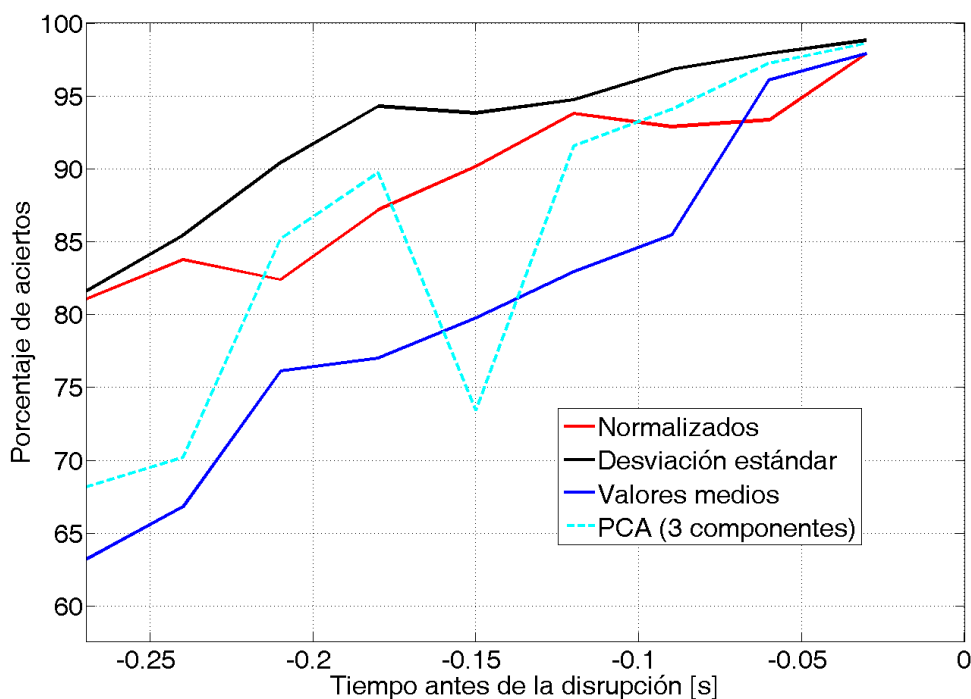


Fig. 5.4 Porcentajes de aciertos de clasificadores SVM entrenados con diferentes características de ventanas disruptivas y no disruptivas.

## 5.4 CONCLUSIONES.

Los diferentes estudios sobre disrupciones realizados en este capítulo concuerdan en que a medida que una descarga se aproxima a la disrupción, ciertos comportamientos inestables del plasma se vuelven más evidentes y por lo tanto más propensos a ser distinguidos. Las técnicas visuales ayudaron a obtener una adecuada extracción de características. Los vectores de características resultantes condensan la información de ventanas temporales de 30 ms vinculadas al instante la disrupción. En una primera instancia se utilizaron estas ventanas como entradas a un sistema de aprendizaje basado en SVM obteniendo altas tasas de aciertos hasta los ~180 ms antes de la disrupción. Este sistema no es aplicable en tiempo real ya que las comparaciones se realizan entre periodos concretos de descargas disruptivas y no disruptivas y no analizan la completa evolución temporal de cada experimento. Este problema, como se verá, es considerado y resuelto en el siguiente capítulo.

---

## **6 Predicción de disrupciones en tiempo real**

---

### **6.1 INTRODUCCIÓN Y ESTADO DEL ARTE.**

En el capítulo anterior se demostró que a través del apropiado procesamiento de las señales es posible condensar en ventanas de 30 ms de duración la información relevante para distinguir si estas se encuentran cercanas en el tiempo a una disrupción. Se determinó la mejor extracción de características que ayudaba a evidenciar el fenómeno y además marcaba claramente que a medida que una descarga se aproxima a la disrupción el comportamiento se vuelve más evidente y por lo tanto más probable de detectar. Sin embargo, el procedimiento anterior no es aplicable al reconocimiento del fenómeno en tiempo real: este no se aplicó a la completa evolución temporal de cada descarga, de principio a fin, como debería realizarse durante la ejecución de un experimento.

En este capítulo se aprovechan los conocimientos adquiridos con el fin de desarrollar un sistema de predicción de disrupciones para JET que, primero, sea aplicable en tiempo real y, segundo, que analice las descargas de principio a fin. Para ello, se emplea la base de datos más extensa que hasta el momento se ha usado en la detección de disrupciones mediante sistemas de aprendizaje.

Con el mismo fin, un reducido número de trabajos basados en técnicas inteligentes ha sido desarrollado previamente como alternativa al sistema de detección

en tiempo real de JET, llamado “Jet Protection System” (JPS) [72]. Entre los más destacados deben mencionarse los de Cannas y col., publicados en los años 2004 y 2007. En el primero de ellos [73] se utiliza una base de datos de 274 descargas. Mediante mapas auto-organizados se realizan agrupamientos en 86 descargas disruptivas para determinar las muestras a emplearse en el entrenamiento del sistema de aprendizaje. Este sistema, basado en redes neuronales, alcanza a detectar un 68% de las interrupciones. En el estudio posterior [74] se menciona que, según la experiencia del equipo investigador, los modelos son capaces de alcanzar altas tasas de acierto únicamente si son probados con descargas pertenecientes al mismo período de las usadas en el entrenamiento. Si en el test se utilizan descargas posteriores, las tasas de acierto caen drásticamente. Proponen como solución un sistema (basado en SVM y acoplado al de predicción) destinado a detectar nuevos comportamientos de manera que estos puedan ser utilizados para un reentrenamiento, evitando el “envejecimiento” del predictor. Tanto la base de datos como el clustering realizado son muy similares a los utilizados en [73]. Se obtienen, sin embargo, porcentajes de reconocimiento más altos: 76%. Allí el 23% de los errores se deben a alarmas perdidas, es decir a descargas disruptivas en las que no se detectó ningún precursor y por lo tanto no se activó la alarma correspondiente.

Fuera de JET, en el Tokamak ASDEX-Upgrade, un estudio de referencia para la predicción del fenómeno es el desarrollado por Pautasso y col. [63]. Se basa en redes neuronales que analizan simultáneamente la evolución de 8 parámetros del plasma y sus derivadas temporales en 99 experimentos. Una vez entrenado, el sistema es probado en un entorno simulado de tiempo real con 500 nuevos pulsos obteniendo un 85% de predicciones correctas con un 1% de alarmas perdidas. En una segunda etapa, las pruebas son realizadas en tiempo real obteniendo un 79% de identificaciones correctas. Otro estudio interesante [64] relaciona dos dispositivos: JET y ASDEX-Upgrade utilizando una base de datos de 185 señales (89 de ASDEX-Upgrade y 96 de JET). En él se entrena un sistema de redes neuronales con descargas de JET y el sistema de predicción es probado en ASDEX-Upgrade con tasas de acierto del 67%. El test inverso logra un porcentaje del 69% de correctas detecciones. En otros Tokamaks también se implementaron sistemas de aprendizaje para la detección de interrupciones. En DIII-D una red neuronal fue entrenada para detectar interrupciones de “alta Beta” [65] intentando establecer umbrales para la activación de alarmas. Se descartaron las descargas en las

que faltaba alguno de los parámetros requeridos por el sistema de detección. La base de datos final estaba constituida por 84 descargas que fueron utilizadas para la etapa de entrenamiento y de pruebas. En el tokamak indio ADITYA [66] un trabajo destinado a detectar interrupciones de límite de densidad utilizaba una base de datos de sólo 23 descargas. Se analizaban únicamente descargas disruptivas y no era aplicable en tiempo real. También con redes neuronales, dos estudios de Yoshino fueron desarrollados en JT-60U. En el primer trabajo [67] el sistema es entrenado en 2 pasos: primero con 12 descargas disruptivas y 6 no disruptivas (paso 1). La información de salida de la red neuronal entrenada se valida y modifica (paso 2) de acuerdo al análisis de un grupo independiente de 12 descargas disruptivas. El modelo final fue probado con 300 descargas disruptivas y 1008 no disruptivas obteniendo tasas de acierto mayores del 80% a 50 ms antes de la interrupción. El segundo trabajo [68] emplea una base de datos de 525 descargas y prueba que un entrenamiento con descargas no disruptivas y un adecuado ajuste del paso 2 se pueden obtener buenos resultados en la detección de interrupciones de mayor dificultad (las descargas de test son del tipo “alta beta”), con una tasa de aciertos del 76%).

Es necesario mencionar los inconvenientes y limitaciones de los trabajos previos:

- Una de las principales limitaciones es la escasa robustez de los sistemas predictores generados. El problema reside en que una vez entrenados, no mantienen su rendimiento en descargas de campañas posteriores a las del entrenamiento. Por consiguiente, la utilidad de los sistemas queda reducida, al no asegurar su capacidad de generalización mediante una tasa de aciertos que permanezca razonablemente estable.
  - En los trabajos previos la única solución propuesta consiste en el continuo reentrenamiento de los sistemas creados, descartando la posibilidad de desarrollar un modelo con la suficiente capacidad de generalización como para funcionar correctamente en campañas posteriores a las de entrenamiento. El continuo reentrenamiento de los sistemas no surge como una solución viable, al ser un proceso lento y a la postre de escasa aplicabilidad.

- En estos estudios la determinación del rendimiento de los sistemas de predicción está limitada por los tiempos que disponen los actuadores para efectuar las acciones de terminación rápida u otras acciones de mitigación. Estas intervenciones de control se basan principalmente en la rápida finalización del calentamiento y en la modificación de la forma y de la corriente del plasma. También se han probado alternativas como insuflar gas o inyectar “*pellets*” con el objetivo de eludir o al menos atenuar los efectos perniciosos de las interrupciones sobre el dispositivo. En general los tipos rápidos de interrupciones son considerados “irreconocibles” y por lo tanto han sido tradicionalmente descartados en el cálculo de las estadísticas finales. Aunque esta clase de interrupciones son escasas, lo cierto es que ignorarlas produce un efecto de maquillaje en los resultados, lo cual debería evitarse.
- Los trabajos desarrollados en JET con sistemas de aprendizaje presentan bases de datos demasiado pequeñas para la obtención de resultados fiables. Asimismo, los estudios en otros Tokamaks (exceptuando [67]) presentan bases de datos considerablemente pequeñas.
- En todos los trabajos citados, las descargas a analizar son seleccionadas cuidadosamente, debiendo cumplir condiciones de estacionalidad y estar limitadas en su máxima corriente. Estas limitaciones son completamente incompatibles si el sistema de predicción quiere ser aplicado en tiempo real.

De lo mencionado anteriormente, quedan claros los objetivos básicos para desarrollar un adecuado sistema de predicción de interrupciones en tiempo real:

- Que sea robusto, asegurando su capacidad de generalización tanto en las campañas en la que es entrenado como en las siguientes.
- Que obtenga tasas de aciertos superiores a las conseguidas hasta el momento en JET, tanto en tiempo real con el JPS como en sistemas diferidos con sistemas de aprendizaje.
- Que sea entrenado y probado mediante una base de datos de gran tamaño, sin sesgar de ninguna manera el tipo de descargas a analizar.

En este capítulo se detallan los pasos seguidos para la consecución de tales objetivos [75 - 77]. Para ello se utilizó una amplísima base de datos, detallada en la Sección 6.2. Con ella se desarrolló un sistema de predicción cuya arquitectura consiste en dos capas. En la primera capa se entrenan una serie de sistemas de clasificación para analizar secuencialmente los experimentos simulando el proceso de tiempo real. Como resultado, cada uno de los clasificadores provee una predicción de la inminencia o no de una interrupción. La determinación de activar o no una alarma dependiendo de los valores de salida provistos por los clasificadores de la primera capa es realizada por otro clasificador también basado en SVM (llamado función de decisión). Esta segunda capa del sistema predictor es un aspecto completamente original y una de las bases de los excelentes resultados obtenidos.

El entrenamiento de la primera capa del sistema de predicción de interrupciones se describe en la Sección 6.3. Se emplean para ello 438 descargas pertenecientes a la campaña experimental C7 de JET (año 2000) y anteriores, las cuales incluyen pulsos desde el 42815 hasta el 57346. Los criterios para adoptar la segunda capa y los detalles de su entrenamiento se describen en la Sección 6.4, mostrándose los resultados en la Sección 6.5.

Las pruebas de la capacidad de generalización del predictor, jamás demostrada hasta el momento en otros trabajos con tasas de éxito comparables, se realiza en la Sección 6.5.2 con descargas pertenecientes a 12 campañas posteriores a las del entrenamiento (desde la campaña C8 hasta la C19, esta última realizada en el año 2007). Los resultados prácticamente no se degradan en sus tasas de aciertos hasta la campaña C14. Después de esta campaña, en el año 2002, el dispositivo fue sometido a cambios significativos, tanto en su estructura como en los diagnósticos encargados de realizar las mediciones de los parámetros del plasma. Aún así las tasas de acierto decaen sólo un 9% (de una media de 88% entre las campañas C8 a C14 a una de 79% de aciertos entre las campañas C15 a C19).

En la Sección 6.6, el rendimiento del predictor es comparado con el del sistema actualmente instalado en JET y por lo tanto la referencia más importante con respecto a la detección de interrupciones. El JPS fue implementado en el año 1991 con el objetivo de efectuar acciones paliativas para anular o al menos mitigar el peligroso efecto de los

abruptos fines de descargas causados por este fenómeno físico<sup>23</sup>. El JPS se basa principalmente en el seguimiento de la señal de “Mode Lock”, la cual mide modos MHD, estableciendo umbrales. Si la señal supera tales umbrales, se activa una alarma. Estos umbrales suelen ser variados manualmente, dependiendo de los propósitos de los experimentos y de los requerimientos de la operación. Su rendimiento dista de ser óptimo. Únicamente detecta un 38% de las interrupciones intencionales (en estos casos el sistema de control vinculado a las alarmas del JPS se encuentra desconectado para que no se inicien las acciones de mitigación). De las interrupciones no intencionales el porcentaje de detección es inferior al 72%. Otro factor a mencionar del JPS es que si éste detecta erróneamente un comportamiento anómalo, incluso no vinculado a una interrupción, se ejecutan las normales acciones de mitigación que a veces son las causantes de las interrupciones en descargas que se estima no iban a interrumpir [78]. El porcentaje de pulsos en los que se cometen estos errores (que podrían ser llamados “falsas alarmas”) no ha sido nunca determinado.

Finalmente, en la Sección 6.7 se discuten los resultados y se plantean las posibles implementaciones del sistema como señal de control en tiempo real en JET.

## 6.2 BASE DE DATOS Y SIMULACIÓN DE TIEMPO REAL.

### 6.2.1 BASE DE DATOS.

Para una mayor fiabilidad de los resultados, se recopiló una extensa base de datos de experimentos de JET (desde el pulso 42815 al 70722) que incluye tanto descargas disruptivas como no disruptivas. Hasta la fecha, es la mayor base de datos de JET destinada al estudio de interrupciones mediante sistemas de aprendizaje en tiempo real. Debido a que JET es el dispositivo de fusión más importante del mundo, los resultados van a tener un impacto relevante en toda la comunidad de fusión.

Un subgrupo de la base de datos (parte sombreada de la Tabla 6.1) fue usado para el entrenamiento y el test del sistema de predicción. Este subgrupo contiene

---

<sup>23</sup> El estudio de la transferencia de las cargas térmicas debidas a las interrupciones a las primeras paredes del dispositivo y la reacción a ellas de los materiales es de gran interés científico. Por ello en cierta cantidad de sesiones experimentales se provocan **intencionalmente** las interrupciones, aunque siempre dentro de configuraciones que aseguren que éstas no causarán daños severos al dispositivo. Las interrupciones **NO intencionales**, en cambio, son aquellas que ocurren inesperadamente durante la rutinaria ejecución de pulsos y deben ser analizadas de manera diferente.



descargas de la campaña C7 y anteriores. Para las pruebas de robustez del predictor, se destinó el resto de descargas, detalladas también en la tabla.

Descargas (Campañas)	Num. de descargas
Disruptivas-entrenamiento (C7 y anteriores)	263
No disruptivas-entrenamiento (C7 y anteriores)	175
Disruptivas-test (C7 y anteriores)	66
No disruptivas-test (C7 y anteriores)	44
Disruptivas test (C1 a C19)	1245
No disruptivas test (C1 a C19)	331
Total disruptivas (C1 a C19)	1574
Total no disruptivas (entrenamiento y test)	550
Total descargas (entrenamiento y test)	2124

*Tabla 6.1 Detalle de la base de datos recopilada para el análisis.*

*En gris se destacan las utilizadas en la etapa de entrenamiento.*

### 6.2.1 SIMULACIÓN DE TIEMPO REAL.

Para la extracción de características, se aprovechó el trabajo desarrollado en el capítulo anterior, y por lo tanto los vectores de características quedan definidos por la desviación estándar del módulo de la transformada de Fourier de las ventanas de 30 ms de la cada señal. El sistema de predicción se diseñó de manera que analice simultáneamente la evolución conjunta de las 13 señales seleccionadas y por lo tanto es necesario que todas y cada una de ellas hayan sido adquiridas durante la ejecución del pulso. En el caso de que alguna de ellas faltase, la descarga no es analizada. El porcentaje de pulsos descartados por este motivo ronda el 8.5%, principalmente debido a ausencias en las mediciones de energía total radiada por el plasma. Sin embargo, si las 13 señales necesarias para este estudio se encuentran disponibles, los experimentos se incluyen en la base de datos para su análisis, incluso si alguna o varias de las

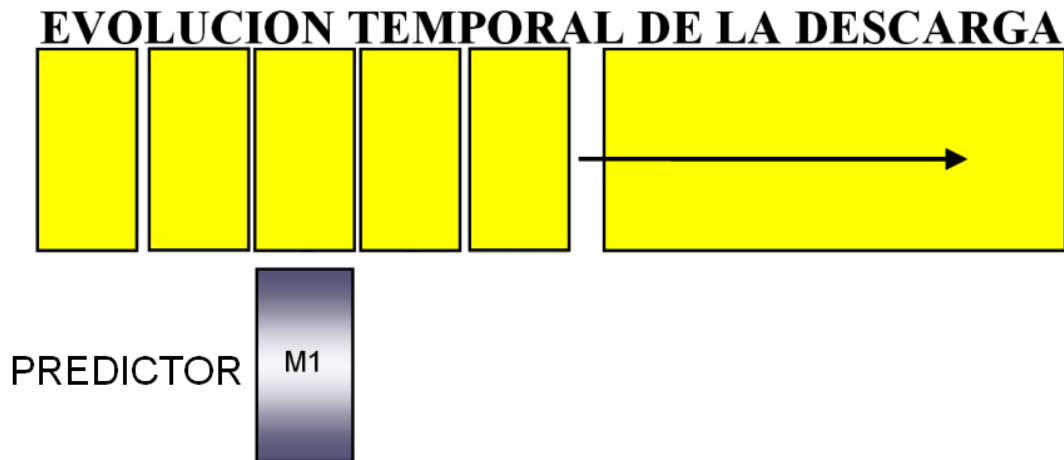
mediciones contienen datos aún no validados. Esto se debe a que el sistema también debe ser capaz de trabajar con algunas mediciones ruidosas o poco fiables. Por motivos de seguridad, en el caso de ser aplicado en tiempo real, el sistema activaría una alarma ante la ausencia de alguna de las señales para detener la descarga.

El objetivo primordial de los estudios sobre el fenómeno es el de evitarlo mediante una detección temprana. Para ello se deben analizar posibles anomalías en las magnitudes del plasma seleccionadas (inestabilidades MHD, incrementos abruptos en la radiación total, cambios inesperados en la corriente o inductancia del plasma, etc.) que puedan indicar comportamientos disruptivos. Además, los sistemas deben ser aplicables en tiempo real. La metodología de las investigaciones anteriores más relevantes sobre el tema (revisadas al comienzo de este capítulo), se basa en el análisis secuencial de las ventanas temporales, tal como se obtendrían durante la ejecución de un experimento.

Para aclarar el concepto, en la Fig. 6.1 se puede observar cómo la total evolución temporal de un pulso puede representarse como la concatenación de vectores de características, cada uno de ellos perteneciente a distintas ventanas temporales, representados por rectángulos. Con un único sistema de clasificación se inspeccionan consecutivamente las ventanas de 30 ms de descarga en busca de posibles precursores. Después de cada análisis el predictor debe tomar una decisión: disparar una alarma o seguir analizando los siguientes 30 milisegundos. Si desde el comienzo hasta el fin del experimento no se decide activar ninguna alarma, esto significa que el sistema no ha reconocido la descarga como disruptiva. En el caso contrario, si un precursor disruptivo es identificado y la alarma es disparada, entonces el sistema ha detectado un comportamiento desencadenante de una interrupción.

La falta de robustez de los clasificadores obtenidos en estudios anteriores (las tasas de acierto decrecen enormemente para campañas de medida diferentes a las de entrenamiento) puede interpretarse de la siguiente manera: los mecanismos que generan el fenómeno disruptivo son la consecuencia de una evolución de inestabilidades demasiado complejas como para ser detectadas satisfactoriamente con un único clasificador, independientemente de su longitud temporal. Por consiguiente, se hace evidente la necesidad de concentrar esfuerzos en el diseño de un sistema capaz de aportar ventajas significativas en cuanto a los resultados. Para ello, y teniendo en cuenta

que la evolución de los precursores pueden presentar escalas temporales de diferente longitud, se implementó la arquitectura completamente innovadora que se detalla en la siguiente Sección.



*Fig. 6.1 La evolución temporal de cada descarga puede simplificarse como la concatenación de los vectores de características de 30 ms.*

*Un sistema predictor puede analizar cada uno de los vectores buscando posibles precursores de interrupciones tal como lo harían en tiempo real.*

### **6.3 ARQUITECTURA DEL CLASIFICADOR Y ENTRENAMIENTO DE LOS MODELOS DE LA PRIMERA CAPA.**

La arquitectura general del sistema de predicción propuesto está compuesta por dos capas. La primera capa, en la que se centra esta Sección, contiene una serie de clasificadores que analizan en paralelo diferentes ventanas de tiempo consecutivas.

El nuevo enfoque, comparado con arquitecturas de un único clasificador pretende:

- Mejorar los resultados obtenidos en los trabajos referenciados en la introducción. Todos esos trabajos se basaban en un único predictor (y por lo tanto constaban de sólo una capa de clasificación).
- Dividir un problema de gran complejidad: este nuevo enfoque se basa en el entrenamiento individual de  $n$  clasificadores para  $n$  intervalos consecutivos de 30 ms.

- Los  $n$  clasificadores resultantes tienen el propósito de analizar en paralelo  $n$  ventanas temporales consecutivas.
- La segunda capa (llamada función de decisión y detallada en la siguiente Sección) evalúa las predicciones de los modelos de la primera y en base a ello activa o no una alarma.

El predictor final puede ser entonces esquematizado en dos capas (Fig. 6.2). En la primera de ellas, la serie de  $n$  modelos analizan  $n$  intervalos consecutivos de descarga obteniendo  $n$  resultados. Cada uno de esos resultados determina si se ha detectado o no comportamiento disruptivo en cada clasificador y lo lejos o lo cerca que se está de la frontera de separación entre los dos regímenes. La segunda capa, detallada en la siguiente Sección y llamada función de decisión (FD), determina mediante el análisis de los  $n$  resultados obtenidos si la alarma debe ser activada o no.

El criterio de entrenamiento de los modelos que conforman la primera capa se basa en diferenciar entre los comportamientos disruptivos o no disruptivos en un instante de tiempo determinado. Dentro del subgrupo de descargas de la base de datos comprendido entre las campañas 7 y anteriores de JET, una proporción balanceada de descargas disruptivas y no disruptivas (263 y 175 respectivamente) se eligieron al azar y fueron destinadas al entrenamiento de los modelos.

La notación utilizada para la identificación de cada modelo puede ser resumida de la siguiente forma:

Sea  $M(i)$  el modelo entrenado para un periodo específico antes de la interrupción, donde  $i=1,2,3,\dots,8$ .  $M(1)$  corresponde al periodo  $[-60, -30]$  ms antes de la interrupción. La ventana correspondiente al intervalo  $[-30, 0]$  antes de la interrupción se descarta para el entrenamiento ya que al terminar en la interrupción no posee información útil para la predicción. Los demás modelos están vinculados a los intervalos  $[-30(i+1), -30i]$  ms antes de la interrupción.

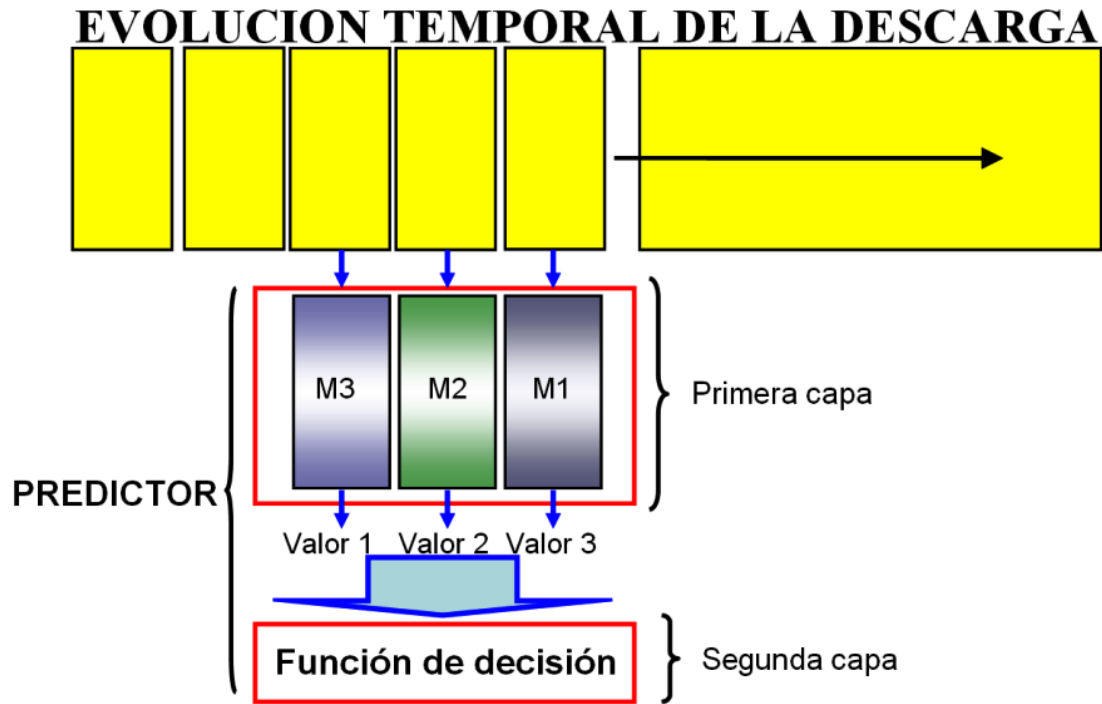


Fig. 6.2 Esquema del análisis en tiempo real mediante una estructura de clasificadores que actúan paralelamente en ventanas de tiempo consecutivas.

En este ejemplo (Secuencia 2) la primera capa está compuesta por 3 modelos.

Como **características disruptivas**, para el modelo M (i):

- Se utilizaron las ventanas correspondientes a los tiempos  $[-30(i+1), -30i]$  ms antes de la interrupción.

Como **características no disruptivas**, para el modelo M (i):

- Se utilizaron todas las ventanas temporales de las descargas no disruptivas.
- También se utilizaron todas las ventanas temporales incluidas en cada descarga 1 segundo antes de la interrupción correspondientes a todos los pulsos disruptivos. Esto se debe a que estadísticamente no hay evidencias de comportamientos disruptivos con una antelación mayor a un segundo. Aún así, como se verá posteriormente, el sistema dispara algunas alarmas (llamadas alarmas prematuras) previas a 1 segundo antes de la interrupción.

En la Fig. 6.3. se muestran las ventanas temporales usadas en el entrenamiento de los modelos de la primera capa. Nótese que para cada uno de ellos la cantidad de

vectores de características no disruptivos (en azul) supera ampliamente a los disruptivos. Este desbalance en la cantidad de vectores se compensa mediante una función interna del programa de entrenamiento [79]. En este software se descartan los vectores de características que aportan la información menos significativa para la detección de interrupciones. Así, los conjuntos finales de entrenamiento tienen una cantidad de ejemplos similar para cada una de las clases considerando los mejores ejemplos de cada uno de los casos.

Para optimizar el método de reconocimiento, se tuvieron en cuenta diferentes combinaciones de los modelos entrenados. A cada una de estas combinaciones se la ha denominado “secuencia”. Al tratarse de un enfoque inédito, la cantidad óptima de modelos a concatenar no podía ser determinada a priori y por lo tanto el óptimo número  $n$  de clasificadores a concatenar era desconocido. Para determinarlo se probaron diferentes combinaciones. Primero se entrenaron 8 modelos<sup>24</sup>. Luego, combinando concatenaciones de los 8 modelos se crearon las siguientes 7 secuencias:

Secuencia 1: M2, M1.

Secuencia 2: M3, M2, M1.

Secuencia 3: M4, M3, M2, M1.

Secuencia 4: M5, M4, M3, M2, M1.

Secuencia 5: M6, M5, M4, M3, M2, M1.

Secuencia 6: M7, M6, M5, M4, M3, M2, M1.

Secuencia 7: M8, M7, M6, M5, M4, M3, M2, M1.

Cada una de las secuencias fue probada, tal como se detallará en las siguientes secciones, siendo una de ellas la que obtenía las tasas de acierto más altas. Para su mejor comprensión, el procedimiento de entrenamiento para la secuencia 2 ha sido resumido en el siguiente pseudo código:

---

<sup>24</sup> El rendimiento de clasificadores entrenados con intervalos previos a 250ms antes de la interrupción es muy bajo debido a que a esos tiempo las señales contienen demasiado poca información como para distinguir con precisión entre una descarga disruptiva o no disruptiva.

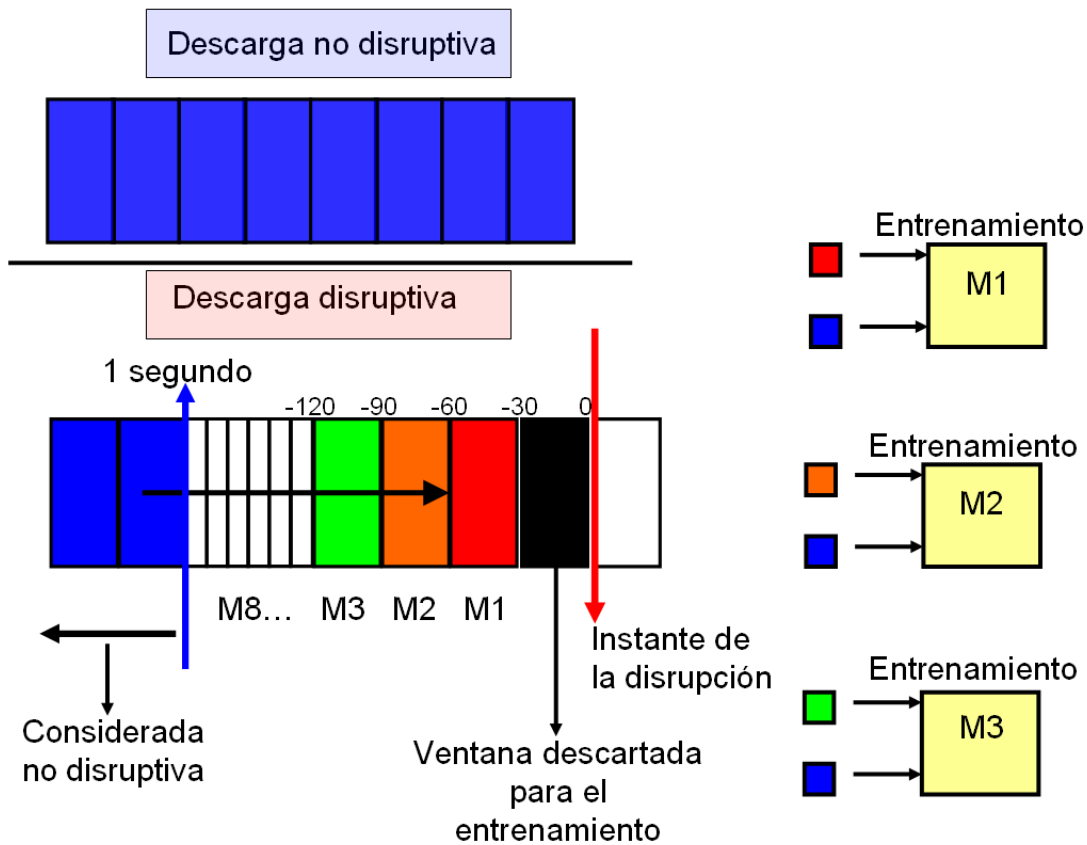


Fig. 6.3 Ejemplo gráfico de las ventanas seleccionadas en las descargas disruptivas y no disruptivas para el entrenamiento de cada uno de los modelos que componen la primera capa del predictor.

MIENTRAS (no se active la condición de final de entrenamiento, detallada en el último párrafo de la Sección 5.4)

PARA (toda la base de datos de entrenamiento)

- 1- El primer modelo de la secuencia (M3) analiza los primeros 30 ms de cada descarga, el segundo modelo los segundos 30 ms y el tercer modelo los terceros 30 ms (ver nuevamente la figura Fig. 6.2).
- 2- Los 3 valores de salida de la secuencia (un valor por clasificador) son analizados por la función de decisión (FD), que consiste en otro clasificador SVM (descrito en detalle en la Sección 5.4). Esta FD determina si la combinación de valores de salida provista por los modelos de la secuencia indican comportamientos disruptivos o no.

SI (se reconoce un comportamiento disruptivo)

Se dispara inmediatamente una alarma y los datos de interés son guardados para un posterior análisis en un formato semejante al representado en la Tabla 6.2. Se analiza la siguiente descarga volviendo al paso 1.

SI NO (no se reconoce un comportamiento disruptivo)

No se dispara alarma alguna ni se guarda ninguna información para posterior análisis. Los siguientes 30 ms del pulso son analizados. Entonces M3 analizará los 30 ms previamente inspeccionados por M2, M2 analizará los 30 ms previamente inspeccionados por M1 y M1 analizará los siguientes 30 ms de descarga. Luego se vuelve al paso 2.

FIN SI

FIN PARA

FIN MIENTRAS

#### **6.4 FUNCIÓN DE DECISIÓN.**

Como se mencionó en la Sección previa, en cada uno de los clasificadores que los  $n$  modelos de una secuencia ejecutan paralelamente sobre tramos de 30 ms consecutivos, se obtienen  $n$  valores de salida. No siempre estos valores coinciden en el carácter disruptivo/no disruptivo de la predicción, ya que algunos pueden detectar un comportamiento anómalo y al mismo tiempo otros no. En consecuencia es necesario implementar una función que evalúe los resultados y que mediante ellos decida si disparar o no una alarma. El desarrollo de esta función de decisión (FD) es crucial para alcanzar la mayor tasa de reconocimiento posible. En la Tabla 6.3. se muestran los valores de los tres modelos de la primera capa que llevaron a la función de decisión a determinar que se debía activar una alarma para tres descargas diferentes. Cada uno de los valores representa la distancia de cada objeto analizado (vector de características de 30 ms) al hiperplano del modelo correspondiente. Distancias positivas corresponden a la clase “disruptiva” y negativas a “NO disruptiva”.

En los tres casos la alarma se activó, correctamente, antes de que ocurriera la interrupción. Nótese que en cada caso, la combinación de valores de los 3 modelos que llevaron a la FD a activar la alarma es muy diferente.



# descarga	Alarma [s]	Disrupción [s]	M3	M2	M1
62999	73.57	73.58	0.57	-0.62	-0.82
63031	44.05	44.14	-1.98	-1.52	0.60
61324	46.27	46.45	-2.21	0.21	0.24

Tabla 6.3 Tres descargas en las que la función de decisión activó una alarma.

#### 6.4.1 ENTRENAMIENTO DE LA FD.

La FD se entrena mediante un proceso iterativo. Únicamente en la primera iteración, un conjunto de condiciones empíricas, llamado Regla Inicial de Decisión (RID) es utilizado para analizar los valores de salida de los modelos y discriminar a través de ellos entre comportamientos disruptivos o no (éste punto será detallado más extendidamente después). En las subsiguientes iteraciones la decisión la tomará la FD. El proceso de entrenamiento se basa en aprovechar los resultados obtenidos en cada iteración previa para refinar la siguiente FD y finaliza al cumplirse una condición de final (comentada al final de la Sección 5.4).

##### 6.4.1.1 La RID.

La RID fue formulada empíricamente teniendo en cuenta la mejor interpretación de los  $n$  valores provistos por los modelos. Para aclarar el concepto se puede decir que si la salida del modelo que analiza el intervalo más cercano en el tiempo a la interrupción no detecta un comportamiento disruptivo (valores mayores a 0), entonces se considera que ese clasificador no ha reconocido ninguna situación anómala. En general, se le asigna una mayor importancia a los valores de ese clasificador, ya que al analizar porciones temporales más cercanas a la interrupción, sus predicciones pueden considerarse como de mayor fiabilidad. La RID se definió independientemente para cada una de las secuencias intentando obtener los mejores resultados posibles. Así, para la secuencia 2 (compuesta por la concatenación de los modelos M3, M2, M1) y cuyos valores de salidas son V3, V2 y V1 respectivamente, la RID empírica es:

$$SI (V3 > -0.8 \ \& \ V2 > -0.4 \ \& \ V1 > 0)$$

Debe dispararse la alarma.

En cambio para la secuencia 1 (M2, M1) las condiciones son:

$SI (V2+V1 > 0)$  o  $SI (V1 > 0.3)$

Debe dispararse la alarma.

A pesar del esfuerzo con la que estos conjuntos de reglas empíricas fue elaborado, resulta obvio que con ellas es improbable que se obtengan los mejores resultados posibles (por ejemplo, para la secuencia 2, la tasa de aciertos no alcanza el 80%).

#### 6.4.1.2 Vectores de característica para el entrenamiento de la FD.

Para entrenar la FD se almacenan los resultados obtenidos en la primera iteración mediante la RID (ver Fig. 6.4). El criterio consiste en guardar los datos relevantes **cada vez que se toma la decisión de disparar una alarma**. Los resultados se almacenan en una tabla que contendrá un número de filas igual al de descargas en las que se activó una alarma. Las predicciones de la primera capa son entonces los vectores de características utilizados para el entrenamiento de la FD.

En la Tabla 6.4 se resumen los 6 posibles tipos de resultados que pueden ser almacenados. La primera columna representa un índice, la segunda el número de descarga, la tercera el tiempo en el que se decide ejecutar la alarma, la cuarta el tiempo en el que ocurre la interrupción (0 si el pulso no es disruptivo), la quinta la diferencia en ms entre el tiempo de la interrupción y el tiempo de la alarma, y las demás los valores de salida de los modelos.

En los ejemplos de la Tabla 6.4 cada índice representa:

1. Un comportamiento disruptivo correctamente detectado 110 ms antes de que esta ocurra.
2. Dos posibles casos de alarmas perdidas (AP):
  - a) Comportamiento anómalo no detectado en una descarga disruptiva.
  - b) Comportamiento anómalo detectado luego de la interrupción.
3. Una alarma prematura (APR): el comportamiento disruptivo es detectado pero demasiado tiempo (más de 1 segundo) antes de la interrupción. En el ejemplo 4339 ms.

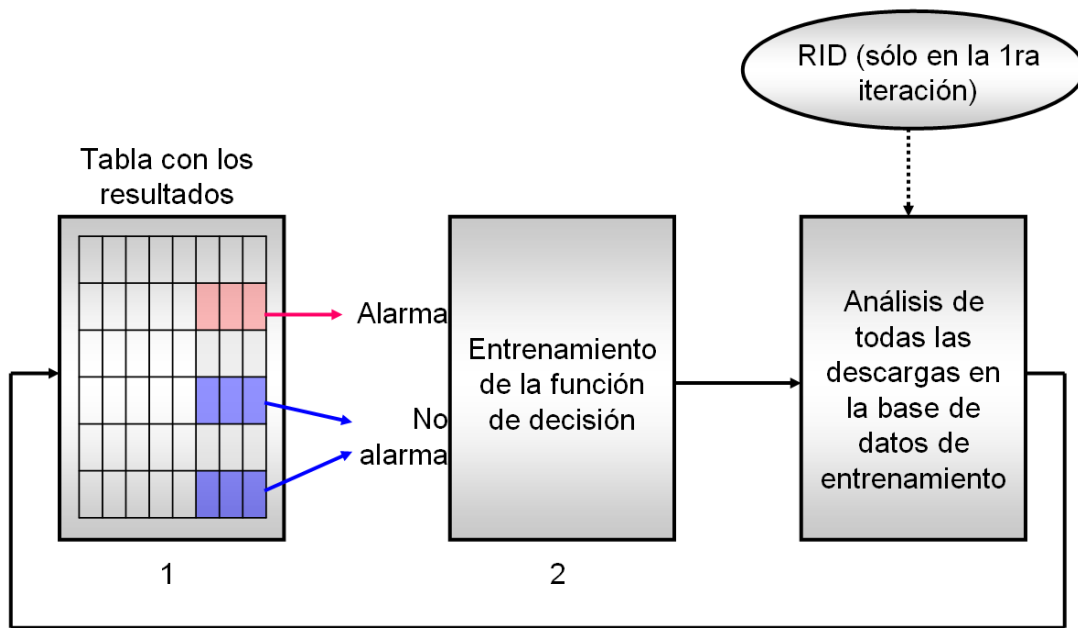


Fig. 6.4 Esquema general de entrenamiento de la función de decisión.

En ausencia de la FD, la RID se utiliza únicamente en la primera iteración y analiza los valores  $V_1$ ,  $V_2$  y  $V_3$  de la primera capa del predictor, creando la primera tabla con los resultados (1). Los valores de la tabla se utilizan para entrenar la función de decisión (2). Finalmente la función de decisión analiza los valores  $V_1$ ,  $V_2$  y  $V_3$  para todas las descargas en las siguientes iteraciones.

4. Un experimento no disruptivo donde correctamente no se ejecutó ninguna alarma.
5. Una falsa alarma (FA): una descarga disruptiva en la cual, incorrectamente, se ejecutó una alarma.

Utilizando los datos almacenados, los valores de salida de los modelos (vectores de características) son introducidos al sistema de aprendizaje SVM para entrenar la FD. Para ello se forman dos grupos, uno conteniendo ejemplos de alarmas activadas correctamente y otro de alarmas activadas incorrectamente.

El primer grupo (**activar la alarma**) contiene los conjuntos de valores de salida de los modelos para el caso 1 (Tabla 6.4, resaltado en rojo), es decir todos los casos en el que el sistema detectó correctamente un comportamiento disruptivo. Con estos ejemplos el sistema es instruido para activar las alarmas en descargas disruptivas.

N°	Pulso	Tiempo de Alarma [s]	Tiempo de Disrupción Time[s]	Margen [ms]	Salida del M3	Salida del M2	Salida del M1
1	56658	63.911	64.021	110	-0.910	0.448	0.0852
2a	54827	0	53.866	0	0	0	0
2b	55253	63.101	62.979	-122	-1.800	-1.391	-1.154
3	53740	45.881	50.22	4339	-0.896	-1	0.445
4	56782	0	0	0	0	0	0
5	52641	40.916	0	-40916	-1.176	-1.202	-1.008

Tabla 6.4. Ejemplos de todos los tipos de resultados que pueden ser almacenados.

Para el entrenamiento de la función de decisión únicamente son útiles los casos resaltados en rojo y azul. Los casos 2<sup>a</sup>, 2b y 4 no aportan valores útiles para que el sistema aprenda cuando es necesario activar o no activar una alarma.

El segundo grupo (**NO activar la alarma**) contiene los conjuntos de valores de salida de los modelos para los casos 3 y 5 de la Tabla 6.4 (azul), con el fin de corregir los errores en la clasificación debido a alarmas falsas y prematuras. Con estos ejemplos el sistema aprende a no activar demasiado pronto las alarmas en descargas disruptivas y a no activarla en descargas no disruptivas.

Las salidas de los modelos M1, M2 y M3 para los casos 2a, 2b y 4 no aportan información útil para el entrenamiento de la FD:

- M1, M2 y M3 para los casos 2a y 4 contienen datos nulos (iguales a 0).
- M1, M2 y M3 para los casos 2b contienen información inútil, ya que corresponde a alarmas activadas demasiado tarde.

Por lo tanto únicamente se entrena la FD, que como se ha dicho antes está basada en SVM, con los vectores de características correspondientes a los valores de M1, M2 y M3 para los casos 1, 3 y 5 de la Tabla 6.4. Una vez concluida la primera iteración y entrenada la primera FD, las decisiones de activar o no la alarma son tomadas por el sistema de aprendizaje, es decir por la FD entrenada. Después de cada

iteración se crean tablas similares a la Tabla 6.4 y los nuevos resultados son añadidos a los previos, teniendo en cuenta no añadir datos repetidos. Esta nueva información es utilizada para entrenar a la siguiente FD. Así, luego de cada bucle, la FD es refinada y mejorada. El proceso de optimización continúa hasta encontrar la condición de finalización. Esta condición es empírica y establece parar el procedimiento tras 5 iteraciones consecutivas en que los porcentajes de acierto del entrenamiento sean menores que el máximo obtenido previamente (ver Fig. 6.5). En general, las máximas tasas se obtuvieron en las primeras 10 iteraciones. Con más iteraciones, y por motivos de sobreajuste, las tasas de acierto de los clasificadores obtenidos disminuía.

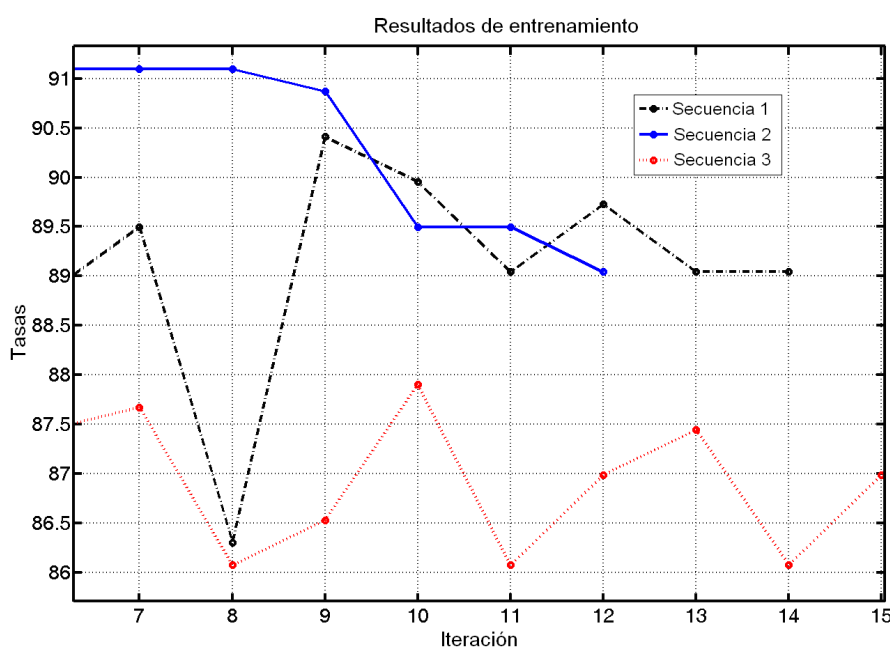


Fig. 6.5 Resultados de las pruebas para tres secuencias.

*El proceso de entrenamiento termina si después de cinco iteraciones no se ha alcanzado un nuevo máximo.*

## 6.5 PRUEBAS DE RENDIMIENTO.

En esta Sección se detallan las pruebas realizadas a los sistemas entrenados. Estas pueden ser divididas en tres diferentes etapas. En la primera, se utiliza un subgrupo de descargas pertenecientes a las mismas campañas en las que el clasificador fue entrenado. Este tipo de test es el que ha sido implementado en trabajos previos y los resultados son útiles para comparar los resultados obtenidos mediante la misma metodología que en otros estudios. La segunda etapa de pruebas consiste en probar el predictor con descargas pertenecientes a campañas posteriores a las de entrenamiento.

Con estas pruebas se pretende demostrar la robustez de sistema. La tercera y última etapa consiste en comparar los resultados con los del sistema actualmente instalado en JET.

#### 6.5.1 PRUEBAS INICIALES.

Para probar el rendimiento del sistema de aprendizaje se realizó una primera serie de pruebas utilizando un subgrupo independiente de descargas. Este tipo de pruebas son las que se realizaron en los estudios anteriores sobre predicción de interrupciones. El objetivo perseguido al realizarlas es el de comparar los resultados en los mismos términos que se obtuvieron en los estudios previos.

Para efectuar las pruebas, tal como ocurría en la etapa de entrenamiento, los pulsos fueron procesados y divididos en ventanas de 30 ms para poder ser analizados por las 7 secuencias de la Sección 6.3. Las pruebas se realizaron seleccionando para cada secuencia la función de decisión que obtenía las mejores tasas de entrenamiento. En la etapa de entrenamiento pudo observarse que los mejores resultados eran obtenidos con la secuencia 2. En la de test los resultados confirmaron esa secuencia como la de mejor rendimiento, tal como se indica en la Tabla 6.5. En la tabla, de izquierda a derecha, cada columna representa las secuencias de clasificadores y los porcentajes de alarmas perdidas (AP), alarmas falsas (AF), alarmas prematuras (APR), total de errores (TE) y tasas totales de aciertos (TA). En la última columna se expresa el tiempo de alarma (tiempo de la interrupción menos tiempo de la alarma) promedio (PROM) en milisegundos para todo el conjunto de test.

#### 6.5.2 PRUEBAS DEL SISTEMA PARA DIFERENTES CAMPAÑAS EXPERIMENTALES.

Para confirmar el potencial de generalización del sistema predictor, este fue probado con descargas pertenecientes a 12 campañas posteriores de JET, desde la campaña C8 hasta la C19. Es necesario aclarar que después de la campaña C14 se realizaron cambios significativos tanto en el dispositivo JET como en algunos sistemas de diagnóstico. Dichas modificaciones, y en especial tres de ellas, afectan directamente al sistema de predicción previamente entrenado. Una de ellas fue la instalación del MKII GB LBSRP (Mark II Gas Box Septum RePlacement), modificación que permitía configuraciones de plasma más elongadas que con las que el predictor había sido entrenado. Otro cambio fue la instalación de Bobinas de Corrección de Error de Campo

	<b>AP</b>	<b>AF</b>	<b>APR</b>	<b>TE</b>	<b>TA</b>	<b>PROM</b>
Predictor 1. [Secuencia 1, DF1]	<b>0</b>	<b>5.4545</b>	<b>5.4545</b>	<b>10.909</b>	<b>89.091</b>	<b>128.01</b>
Predictor 2. [Secuencia 2, DF2]	<b>.909</b>	<b>4.5455</b>	<b>1.8182</b>	<b>7.2727</b>	<b>92.727</b>	<b>146.23</b>
Predictor 3. [Secuencia 3, DF3]	<b>0</b>	<b>4.5455</b>	<b>5.4545</b>	<b>10</b>	<b>90</b>	<b>132.23</b>
Predictor 4. [Secuencia 4, DF4]	<b>0</b>	<b>5.4545</b>	<b>9.0909</b>	<b>14.545</b>	<b>85.455</b>	<b>128.84</b>
Predictor 5. [Secuencia 5, DF5]	<b>0</b>	<b>4.5455</b>	<b>7.2727</b>	<b>11.818</b>	<b>88.182</b>	<b>136.45</b>
Predictor 6. [Secuencia 6, DF6]	<b>0</b>	<b>5.4545</b>	<b>5.4545</b>	<b>10.909</b>	<b>89.091</b>	<b>136.23</b>
Predictor 7. [Secuencia 7, DF7]	<b>0</b>	<b>5.4545</b>	<b>6.3636</b>	<b>11.818</b>	<b>88.182</b>	<b>121.01</b>

*Tabla 6.5 Resultados de las pruebas para todas las secuencias.*

*A pesar de que la secuencia 2 tiene un porcentaje no nulo de alarmas perdidas, el total de aciertos es el mayor y por lo tanto puede considerarse como el mejor sistema predictor.*

(BCEC) [80], que pueden afectar a las mediciones de la señal de “Mode Lock” empleada por el predictor. Además, se reemplazó el sistema de bolometría (del cual el predictor utiliza la señal de potencia radiada). El nuevo sistema adquiere formas de onda con diferencias significativas respecto a la del diagnóstico anterior (en el sistema más reciente pueden observarse picos en las señales debidos a ELMs).

Los resultados del sistema bicapa detector de interrupciones en JET han sido resumidos en la Fig. 6.6. En ellos se detallan únicamente las tasas totales de aciertos y las falsas alarmas (las alarmas prematuras y perdidas se mantienen constantes y por lo tanto no aportan información útil). Cabe aclarar que la suma de estos porcentajes no alcanza el 100%. Esto se debe a que no se ha representado la evolución de los errores debido a alarmas perdidas o alarmas prematuras, ya que estos valores se mantienen estables (la suma de ambos no supera el 11%).

En la Fig. 6.6, la zona pintada de verde indica la última campaña en la que el modelo fue entrenado y testeado (Sección 6.5).

Sin sombrear, en el centro de la gráfica, se muestran los porcentajes calculados para 376 descargas (la mitad de ellas disruptivas). La variación de la tasa de aciertos es menor del 11.6% (entre 94.12% y 82.35%). Debido al reducido número de descargas

disponibles de la corta campaña 9, ésta fue incluida en las estadísticas con la campaña 8 (fueron empleadas en total 100 descargas, 30 de ellas disruptivas). En este intervalo (entre las campañas C8 y C14) la mínima tasa de aciertos se obtiene en la campaña C11, llamada “Trace Tritium Campaign”. Por el uso de tritio como combustible, se impusieron restricciones en la operación del dispositivo para evitar aquellas configuraciones en las que las interrupciones suelen ocurrir con mayor frecuencia. Consecuentemente el número de interrupciones para esta campaña no sólo fue menor, sino que se evitaron las interrupciones más típicas y la mayoría de las ocurridas fueron de los tipos más rápidos y por lo tanto las más difíciles de detectar.

Finalmente la parte de la gráfica sombreada de azul muestra los resultados obtenidos entre las campañas C15 y C19 (246 descargas, la mitad de ellas disruptivas). Como se mencionó, las causas de las menores tasas de acierto están vinculadas a las modificaciones realizadas en la máquina. De esta forma se puede observar, de forma realista, cómo pueden influir tales cambios en los resultados del predictor. Los peores porcentajes se obtienen en la campaña C19, donde se evidencia un incremento considerable de las falsas alarmas. Alrededor de un 47% de esas alarmas, erróneamente activadas por el predictor, son causadas por el uso de las BCEC con el consecuente efecto sobre la señal de “Lock mode”. Otro ~44% de las falsas alarmas fueron activadas por la identificación de picos de radiación que únicamente podían ser medidos por el nuevo bolómetro.

Vale la pena mencionar que la robustez demostrada por el predictor se mantiene durante un amplio rango de campañas, especialmente entre la C8 y la C14, ninguna de ellas pertenecientes al período en el que el predictor fue entrenado. Estos resultados no habían sido obtenidos hasta el presente y por lo tanto debe asignárseles una gran importancia.

Incluso después de la campaña C14 las tasas son altamente satisfactorias ya que las causas de la mayoría de los fallos se deben a falsas alarmas debidas a cambios considerables realizados en el dispositivo.



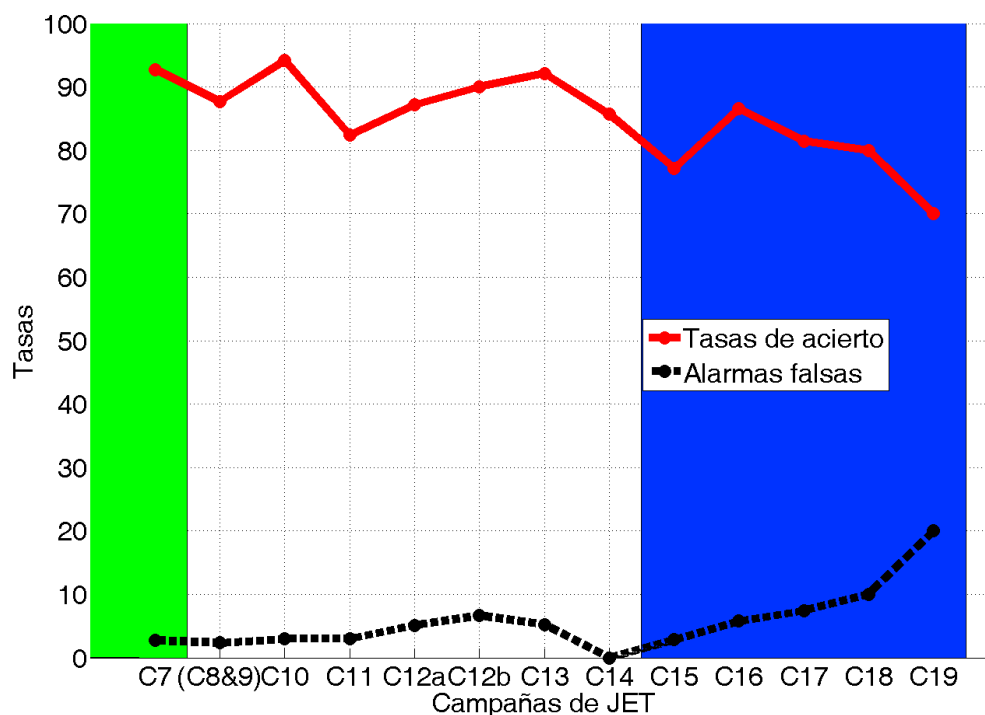


Fig. 6.6 Tasas de acierto del modelo previamente entrenado.

Tres diferentes periodos son representados. A la izquierda se resalta el rango de campañas (C7 y anteriores) donde el predictor fue entrenado. En la parte sin sombrear se representan los resultados obtenidos por el predictor sin reentrenamiento para el período comprendido entre las campañas C8 y C14. Finalmente en azul se representan los porcentajes del predictor para las campañas C15 a C19.

## 6.6 COMPARACIÓN DEL RENDIMIENTO DEL SISTEMA PREDICTOR DESARROLLADO CON EL SISTEMA JPS.

El rendimiento del predictor desarrollado ha sido comparado con el sistema JPS. Para contrastar apropiadamente los resultados obtenidos por cada sistema fue necesario tener en cuenta ciertas consideraciones preliminares. El sistema JPS interviene durante la ejecución de un experimento cada vez que identifica algún comportamiento disruptivo, desencadenando una serie de acciones de mitigación o apagado rápido. Tanto si el sistema está detectando correctamente un comportamiento disruptivo como si el JPS erróneamente ha disparado una falsa alarma, se interviene inmediatamente sobre la descarga. No existen por lo tanto estadísticas sobre falsas alarmas, ya que cada alarma activada conlleva al apagado del experimento. Es más, en algunos casos las acciones de atenuación y finalización rápida del pulso provocadas por alarmas erróneas pueden llevar a que descargas no disruptivas terminen en una interrupción [78]. Las estadísticas entonces serán únicamente calculadas para pulsos disruptivos. También por motivos comparativos, las tasas de aciertos no son calculadas como en las secciones

anteriores, en las que se desglosaban los errores en alarmas falsas, prematuras y perdidas, sino de la misma forma en la que suelen calcularse para el sistema de protección de JET.

Las curvas de la Fig. 6.7 representan los porcentajes acumulados de las interrupciones detectadas (eje de ordenadas) para diferentes “tiempos de alarma” (abscisas en escala logarítmica). Estos tiempos de alarma (instante de la interrupción menos instante en el que el sistema ejecutó la alarma) son de gran importancia ya que definen el margen temporal que los actuadores disponen para realizar acciones de apagado o mitigación. Para mejorar la comprensibilidad de los resultados obtenidos, se calcularon estadísticas independientes entre las interrupciones no intencionales (las que ocurren durante condiciones normales de operación) e intencionales (aquellas que son provocadas intencionalmente para estudiar la física del fenómeno). Las curvas de la Fig. 6.7.a representan las estadísticas sobre todas las interrupciones intencionales ocurridas desde la campaña C1 hasta la C19 (347 interrupciones). Las mismas estadísticas, pero para las interrupciones no intencionales (136 interrupciones), se calculan y se muestran en la Fig. 6.7.b. En la Fig. 6.7.c los porcentajes son computados para las 483 interrupciones (es decir para las todas las intencionales y no intencionales). Nótese en esta última representación gráfica, que el sistema desarrollado, con márgenes de tiempo de 30 ms detecta el 77% de las interrupciones mientras que el sistema JPS solamente un 48%. Los porcentajes de detección son mayores hasta los 200 ms aproximadamente. Estos dos tiempos de alarma son especialmente relevantes en JET, ya que en este dispositivo los tiempos típicos para efectuar las acciones de control más comunes, como la modificación de la forma y de la corriente del plasma, requieren 30 ms y 200 ms respectivamente. Es necesario mencionar, a favor del JPS, que este sistema computa las estadísticas sobre absolutamente todas las descargas, incluso si en algunas de ellas se omite la adquisición de alguna de las señales, mientras el predictor desarrollado descartó un 8.5% de los experimentos debido a la falta de algunas de las señales necesarias para su correcto funcionamiento. Es destacable que a pesar de que el sistema de predicción fue entrenado para no disparar alarmas con más de 1 segundo de antelación, tal como puede verse en las figuras anteriores existen excepciones en las que sí se activan alarmas, probablemente debidas a cambios significativos en los parámetros del plasma.

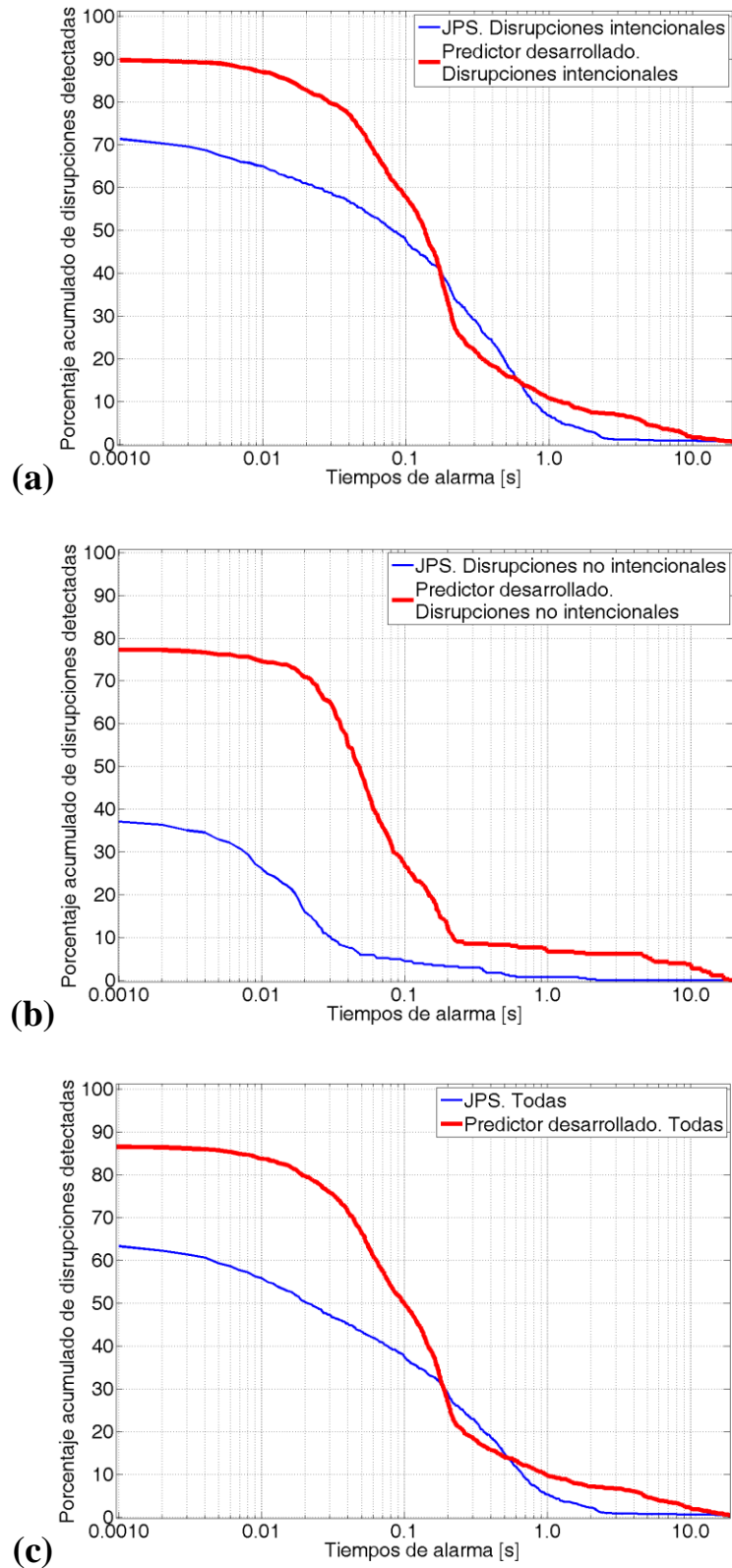


Fig. 6.7 Porcentajes acumulados de interrupciones detectadas para diferentes tiempos de alarma.  
 a) Estadísticas calculadas para todas las descargas disruptivas no intencionales desde la campaña C1 hasta la C19. b) Las mismas estadísticas para descargas cuyas interrupciones fueron provocadas intencionalmente. c) Estadísticas conjuntas (incluyen a) y b).

## 6.7 CONCLUSIONES.

En este capítulo se concatenaron varios modelos específicamente entrenados con SVM para obtener 7 secuencias. Con ellas, y con una función de decisión también basada en SVM, se obtuvieron altas tasas de reconocimiento. El predictor final consiste en una novedosa arquitectura bi-capa capaz de trabajar en tiempo real. La secuencia que demostró una mayor tasa de aciertos fue la secuencia 2. El sistema fue testado en tres etapas diferentes.

En primer lugar, para las campañas 7 y anteriores, es decir las mismas en las que el predictor fue entrenado, la tasa de aciertos en el grupo reservado para pruebas fue del 92.73%, con menos de un 0.91% de alarmas perdidas. Este porcentaje de alarmas perdidas representa 1 error en 110 descargas. La arquitectura general del sistema y el entrenamiento minimizan el número de alarmas perdidas pero también acotan los otros posibles tipos de errores (la suma de alarmas falsas y alarmas prematuras queda reducida a menos del 6.4%).

En segundo lugar se comprobó la robustez del sistema con un gran número de descargas correspondientes a campañas más recientes de JET que con las que había sido entrenado. Este tipo de pruebas tiene una alta importancia ya que la falta de robustez es una de las mayores debilidades de los modelos desarrollados anteriormente. Los resultados muestran claramente que las tasas de acierto se mantienen altas hasta que cambios estructurales de importancia son efectuados sobre la máquina. Después de dichos cambios las tasas de acierto continúan siendo satisfactorias, a pesar de caer un 9%.

Finalmente el rendimiento del sistema es comparado con el sistema implementado en JET desde hace muchos años, el JPS. Las tasas de reconocimiento, para tiempos de alarma hasta los ~200 ms son considerablemente más altas para el método basado en la arquitectura bicapa. Estos resultados son de gran interés ya que el tiempo de respuesta del sistema de mitigación que controla la forma del plasma tiene un tiempo de ejecución típico de entre 30 ms y 200 ms.

Los resultados fueron probados con la más amplia base de datos que haya sido usada para sistemas automáticos de reconocimiento de interrupciones en tiempo real. En esta investigación se le asignó la máxima prioridad a lograr la mayor tasa de aciertos

posible minimizando las alarmas perdidas y teniendo en cuenta en segundo lugar el margen de tiempo de acción.

La aplicación en tiempo real del sistema predictor ya está siendo llevada a cabo en JET y se espera su operación rutinaria en las próximas campañas.



---

## 7 Conclusiones

---

### 7.1 CONCLUSIONES.

En esta Tesis se desarrollaron diferentes sistemas de minería de datos a señales provenientes de dispositivos destinados al estudio de la fusión nuclear. Estos sistemas se aplicaron a diferentes bases de datos, en particular las de TJ-II y JET (cuyo repositorio contiene más de 40 Tbytes de información).

Cualquier evento físico en el plasma genera normalmente los mismos patrones de señal. Por lo tanto, es posible establecer búsquedas de formas que representen fenomenologías específicas. Los algoritmos de búsqueda desarrollados en esta tesis permiten este tipo de búsquedas y deben considerarse herramientas particularmente necesarias en las masivas bases de datos de fusión. Su aplicación, además de sencilla es eficiente.

Tan relevante como las búsquedas de patrones en las señales es la interpretación de la información hallada. Para encontrar la interrelación de magnitudes del plasma en el caso de las transiciones L-H se utilizó SVM. Este sistema de aprendizaje, para clasificación, construye un hiperplano que detalla la relación entre los parámetros con

los que el sistema fue entrenado. La ecuación del hiperplano hallada representa un modelo lineal con el que puede determinarse el umbral en el que se producen las transiciones. La metodología posibilita la extracción de ecuaciones que ayudan a la comprensión de fenómenos físicos.

Un meticuloso estudio sobre interrupciones fue llevado a cabo en JET con el fin de poder prevenirlas durante la operación del dispositivo. Se estudió la adecuada extracción de características para finalmente entrenar un sistema de predicción aplicable en tiempo real. El predictor obtenido, completamente original en su arquitectura, logra alcanzar las tasas de acierto más altas hasta el momento. Un factor aún más importante es su robustez. Al respecto, es relevante recordar que en fusión los periodos de investigación se estructuran en campañas experimentales. En los trabajos desarrollados anteriormente, tanto en JET como en ASDEX-Upgrade, los sistemas se entrenaban con un conjunto de descargas pertenecientes a una campaña. Estos predictores jamás demostraron capacidad para obtener resultados satisfactorios cuando eran probados con experimentos correspondientes a campañas diferentes. El sistema desarrollado en esta tesis, al contrario, es capaz de mantener razonablemente estables las tasas de aciertos con descargas hasta 7 años (y 12 campañas) posteriores a las del conjunto de entrenamiento. Incluso, estos buenos resultados se conservan después de que el dispositivo hubiese sufrido cambios estructurales importantes. La relevancia de este predictor es muy significativa, ya que mejora sensiblemente la detección oportuna de uno de los principales inconvenientes para la obtención de energía termonuclear en Tokamaks. El sistema está siendo instalado actualmente para su uso en tiempo real en JET.

## 7.2 TRABAJOS FUTUROS.

El sistema de predicción de interrupciones puede ser mejorado. En particular se pretende trabajar en el futuro en tres aspectos particulares:

- Aportar valores de credibilidad y confianza para cada predicción. Al respecto, nuevas teorías estadísticas desarrolladas por Vovk y colegas [83] pueden ser aplicadas.



- Detectar el tipo de interrupción: si al detectar un comportamiento disruptivo se es capaz además de determinar el tipo de interrupción se podrían aplicar medidas de mitigación diferentes para cada caso.
- Profundizar en la selección del conjunto de señales y en el estudio del adecuado proceso de extracción de características para el entrenamiento del sistema. Actualmente se están desarrollando metodologías que pueden ser útiles para tal propósito.



---

## 8 Bibliografía

A continuación se presenta la bibliografía utilizada.

- [1] C.M. Braams y P.E. Stott. “Nuclear Fusion. Half a Century of Magnetic Confinement Fusion Research”. *Plasma Phys. Control. Fusion* 44 1767. 2002
- [2] J. Sheffield. “The physics of magnetic fusion-reactors”. *Rev. Mod. Phys.* 66, 1015. 1994,
- [3] J. R. Reitz, F. J. Milford y R. W. Christy, “*Fundamentos de la Teoría Electromagnética*”. Addison-Wesley Iberoamericana, Wilmington. 1996.
- [4] J. D. Lawson, “*Some criteria for a power producing thermonuclear reactor*”, *Proc. Phys. Soc. B*, 70, 6 (1957).
- [5] A. R. Bell, “*Laser Produced Plasmas*”, *Plasma Physics: An Introductory Course*. Cambridge University Press, Cambridge (1993).
- [6] McCracken G. and Stott P., *Fusion: The Energy of the Universe*, Elsevier Academic Press, 2005.

[7] M. Wakatani, “Stellarator and Heliotron Devices”. Oxford University Press, Oxford 1998.

[8] S. Yoshikawa, “Design of a helical-axis stellarator”, Nucl. Fusion, 23, 667 (1983)., J. H. Harris, J. L. Cantrell, T. C. Hender, B. A. Carreras y R. N. Morris, “A flexible heliac configuration”, Nucl. Fusion, 25, 623. 1985.

[9] T. C. Hender, J. F. Lyon, J. L. Cantrell, J. A. Fábregas, J. H. Harris, J. Guasp, B. A. Carreras, A. López-Fraguas, V. E. Lynch y A. P. Navarro, “Studies of a flexible heliac configuration”, Fusion Technol., 13, 521. 1988.

[10] C. Alejaldre, J. Alonso, L. Almoguera, E. Ascasíbar, A. Baciero, R. Balbín, M. Blaumoser, J. Botija, B. Brañas, E. de la Cal, A. Cappa, R. Carrasco, F. Castejón, J. R. Cepero, C. Cremy, J. Doncel, C. Dulya, T. Estrada, A. Fernández, M. Francés, C. Fuentes, A. García, I. García-Cortés, J. Guasp, J. Herranz, C. Hidalgo, J. A. Jiménez, I. Kirpichev, V. Krivenski, I. Labrador, F. Lapayese, K. Likin, M. Liniers, A. López-Fraguas, A. López-Sánchez, E. de la Luna, R. Martín, A. Martínez, M. Medrano, P. Méndez, K. J. McCarthy, F. Medina, B. van Milligen, M. Ochando, L. Palacios, P. Pastor, M. A. Pedrosa, A. de la Peña, A. Portas, J. Qin, L. Rodríguez-Rodrigo, A. Salas, E. Sánchez, J. Sánchez, F. L. Tabarés, D. Tafalla, V. Tribaldos, J. Vega, B. Zurro, D. Akulina, O. I. Fedyanin, S. Grebenschikov, N. Kharchev, A. Meshcheryakov, R. Barth, G. van Dijk, H. van der Meiden y S. Petrov, “*First plasmas in the TJ-II flexible heliac*”, Plasma Phys. Control. Fusion, 41, A539. 1999.

[11] R. A. Cairns, “*Radio-frequency Plasma Heating*”, Plasma Physics: An Introductory Course. Cambridge University Press, Cambridge. 1993.

[12] F. L. Tabarés, D. Tafalla, R. Balbín, B. Brañas, T. Estrada, I. García-Cortés, F. Medina y M. A. Ochando, “*Impact of wall conditioning and gas fuelling on the enhanced confinement modes in TJ-II*”, J. Nucl. Mater., 313-316, 839. 2003.

[13] K. Miyamoto, “*Plasma Physics for Nuclear Fusion*”, MIT Press, Cambridge. 1980.

[14] <http://www.jet.efda.org/>

- [15] J Ongena. “JET’s contribution to fusion science and ITER”. Phys. Scr. T123 14–23. (2006)
- [16] Pearson, K. (1901). “On Lines and Planes of Closest Fit to Systems of Points in Space” (PDF). Philosophical Magazine 2 (6): 559–572.  
<http://stat.smmu.edu.cn/history/pearson1901.pdf>.
- [17] Gentle, J. E. “Elements of Computacional Satatistics”. New Cork: Springer-Verlag. 2002.
- [18] Strang, Gilbert. “Introduction to Linear Algebra”. Wellesley, MA: Wellesley-Cambridge Press. 1993.
- [19] Cox, Trevor and Cox, Michael. “*Multidimensional Scaling, 2<sup>nd</sup> Edition*”. Boca Raton: Chapman & Hall/CRC. 2001.
- [20] W.L.Martinez and A.R.Martinez “Exploratory data Analysis with Matlab” Chapman&Hall. 2005.
- [21] Asimov, Daniel. “The grand tour: A tool for viewing multidimensional data”. SIAM Journal of Scientific and Statistical Computing. 6 :128-143. 1985.
- [22] Haykin S. “Neural Networks: A Comprehensive Foundation”. New Jersey:Prentice Hall. 1999.
- [23] Bishop, C. Svensèn, M. Williams, C. “*The Generative Topographic Mapping*”, Neural Computation 10, No. 1, 215-234. 1998.
- [24] J. B. MacQueen. “*Some Methods for classification and Analysis of Multivariate Observations*”, Proceedings of 5-th Berkeley Symposium on Mathematical Statistics and Probability, Berkeley, University of California Press, 1:281-297). 1967.

[25] Bradley, P.S. y Fayyad, U.M. “*Refining initial points for k-means clustering*”. In J. Shavlik, editor, Proceedings of the Fifteenth International Conference on Machine Learning (ICML '98), pages 91--99, San Francisco, CA. Morgan Kaufmann. 1998.

[26] Roger J. Lewis. “*An Introduction to Classification and Regression Tree (CART) Analysis*”. Annual Meeting of the Society for Academic Emergency Medicine in San Francisco, California. 2000.

[27] Corinna Cortes and V. Vapnik, “Support-Vector Networks”, Machine Learning, 20. 1995.

[28] V. Vapnik. “The Nature of Statistical Learning Theory”. Springer, Second Edition. 2000.

[29] <http://www2.sims.berkeley.edu/research/projects/how-much-info/>

[30] J. Vega, G.A. Rattá, P. Castro, A. Murari. “Development of Learning Systems with Data Tours Techniques for Fusion Databases”. Proceedings of the 8th International FLINS Conference, Madrid, España. pps 103-109. 2008.

[31] H. Nakanishi et al. “Search and retrieval method of similar plasma waveforms”. Fusion Engineering and Design 71. 189–193. (2004)

[32] J. Vega “Intelligent methods for data retrieval in fusion databases”. Fusion Engineering and Design 83. 382–386. 2008.

[33] J. Vega et al. “Data mining technique for fast retrieval of similar waveforms in Fusion massive databases”. Fusion Engineering and Design, vol. 83, pp. 132-139. 2008.

[34] S. Dormido-Canto, G. Farias, J. Vega, R. Dormido, J. Sánchez, M. Santos, J. A. Martín, G. Pajares.. “Search and retrieval of plasma wave forms: Structural pattern recognition approach”. Review of Scientific Instruments 77, 10F514. 2006.

- [35] G.A. Rattá, J. Vega, A. Pereira, A. Portas, E. de la Luna, S. Dormido-Canto, G. Farias, R. Dormido, J. Sánchez, N. Duro, H. Vargas, M. Santos, G. Pajares, A. Murari. “First applications of structural pattern recognition methods to the investigation of specific physical phenomena at JET”. *Fusion Engineering and Design* Volume 83, Issues 2-3. Pages 467-470. April 2008.
- [36] S. Dormido-Canto, G. Farias, R. Dormido, J. Vega, J. Sánchez, N. Duro, H. Vargas, G. Rattá, A. Pereira and A. Portas. “Structural pattern recognition methods based on string comparison for fusion databases”. *Fusion Engineering and Design*. Volume 83, Issues 2-3. 2007.
- [37] A. Pereira, J. Vega, A. Portas, R. Castro, A. Murari. “Optimized search strategies to improve structural pattern recognition techniques”. *Proceedings of the 8th International FLINS Conference, Madrid, España*. pps 405-411. 2008.
- [38] J. Vega, A. Murari, G. A. Rattá, P. Castro, A. Pereira, A. Portas and JET-EFDA Contributors. “Structural pattern recognition techniques for data retrieval in massive fusion databases”. *Burning Plasma Diagnostics. AIP Conference Proceedings* (ISBN 978-0-7354-0507-3). 988 481-484. 2008.
- [39] J. Vega, A. Murari, A. Pereira, A. Portas, G. A. Rattá, R. Castro and JET-EFDA Contributors. “Overview of intelligent data retrieval methods for waveforms and images in massive fusion databases”. *Fusion Engineering and Design* 84 1916–1919. 2009.
- [40] J. Vega, G. A. Rattá, A. Murari, P. Castro, S. Dormido-Canto, R. Dormido, G. Farias, A. Pereira, A. Portas, E. De la Luna, I. Pastor, J. Sánchez, N. Duro, R. Castro, M. Santos, H. Vargas. “Recent results on structural pattern recognition for Fusion massive databases”. *Proc. of the IEEE International Symposium on Intelligent Signal Processing*. ISBN: 1-4244-0829-6 (2007) 949-954.
- [41] E. de la Luna et al. “Electron cyclotron emission radiometer upgrade on the JET Tokamak”. *Rev. Sci. Instrum.* 75 (10). 3831-3833. 2004.
- [42] J. Wesson, “Tokamaks”. Oxford: Clarendon Press Oxford. Third Edition. 2004.

- [43] Wagner, F. et al., Physical Review Letters 49, 1408. 1982.
- [44] J.Vega, A. Murari, G. Vagliasindi, G. A. Rattá and JET-EFDA Contributors. “Automated estimation of L/H transition times at JET by combining Bayesian statistics and Support Vector Machines”. Nucl. Fusion 49 (2009) 085023 (11pp). 2009.
- [45] M. Ruiz, J. Vega, E. Barrera, J. González, A. Murari, R. Meléndez, G.Rattá, S. González and JET-EFDA Contributors. “Test-bed of a real time detection system of L-H & H-L plasma transitions implemented with ITMS platform”. (To be submitted to Fusion Eng. & Design) 2009.
- [46] <http://efdsql.ipp.mpg.de/HmodePublic/>
- [47] "The Spider - A machine learning in Matlab".  
<http://www.kyb.tuebingen.mpg.de/bs/people/spider/main.html> Max-Planck Institute for biological Cybernetics, Tuebingen, Alemania.
- [48] Chankin A. V.. Plas. Phys. Cont. Fus. 39 1059. 1997.
- [49] Chankin A. V., Saibene G. Plas. Phys. Cont. Fus. 41 913. 1999.
- [50] Kerner W. et al. Contrib. Plas. Phys. 38 118. 1998.
- [51] Rogister, A. L. Plas. Phys. Cont. Fus. 36 A219. 1994.
- [52] Scott B et al. Fusion Energy (Proc. 16th Int. Conf., Montreal) vol 2 (Vienna: IAEA) 649. 1996.
- [53] Shaing K C, Crume Jr. E C. Phys. Rev. Lett. 63 2369. 1989.
- [54] A.Murari, J.Vega, D.Mazon, N.Martin, G.Rattá, G.Vagliasindi. “Machine Learning Methods for Data Driven Theory in the Physical Sciences with Applications to Confinement Regime Identification in Nuclear Fusion”. Enviado a Reports on Progress in Physics.



- [55] “Chapter 3: MHD stability, operational limits and disruptions”. ITER Physics Expert Group on Disruptions, Plasma Control, and MHD et al. Nuclear Fusion, Vol. 39, No. 12. 1999.
- [56] B. Cannas, et al. “Disruptions forecasting at JET using Neural Networks”. Nucl. Fusion. 44 68-76. 2004
- [57] B. Cannas, et al. “Disruption prediction at ASDEX Upgrade using neural Networks”. 33rd EPS Conference on Plasma Physics. Contributed Papers, (Eds.) F. De Marco, G. Flad. ECA 30 I. European Physical Society, Geneva, P-2.143. 2006.
- [58] A. Murari, et al. “Prototype of an adaptive disruption predictor for JET based on fuzzy logic and regression trees”. Nucl. Fusion 48 035010. 2008.
- [59] G.A. Rattá, J. Vega, A. Murari and M. Johnson. “Feature extraction for improved disruption prediction analysis at JET”. Rev. Sci. Instrum. 79, 10F328. 2008.
- [60] F.C. Schuller. “Disruption in Tokamaks”. Plasma Phys. Control. Fusion. 37 A135–62. 1995.
- [61] Alexei Savtchkov. “Mitigation of disruptions in a Tokamak by means of large gas injection”. Jülich. 2003.
- [62] A.W. Morris. “MHD instability control, disruptions, and error fields in Tokamaks”. PPCF Vol 34. N 13. Pps 1871-1879. 1992
- [63] Pautasso G. et al. “On-line prediction and mitigation of disruptions in ASDEX Upgrade”. Nucl. Fusion 42 100–08. 2002.
- [64] Windsor C.G. et al 2005 Nucl. Fusion 45 337.
- [65] Wroblewski D., Jahns G.L. and Leuer J.A. 1997 Nucl. Fusion 37 725.
- [66] Sengupta A. and Ranjan P. 2001 Nucl. Fusion 41 487.

- [67] Yoshino R. 2003 Nucl. Fusion 43 1771.
- [68] Yoshino R. 2005 Nucl. Fusion 45 1232.
- [69] A. Murari, J. Vega, G.A. Rattá, G. Vagliasindi, M.F. Johnson, S.H. Hong. “Unbiased and Non-Supervised Learning Methods for Disruption Prediction at JET”. Nucl. Fusion 49. 055028. 2009.
- [70] A. Murari, J. Vega, J. A. Alonso, E. de la Luna, J. Farthing, C. Hidalgo, G. A. Rattá, J. Svensson, G. Vagliasindi and JET EFDA Contributors. “How to extract information and knowledge from fusion massive databases”. Burning Plasma Diagnostics. AIP Conference Proceedings (ISBN 978-0-7354-0507-3). 988 457-470. 2008.
- [71] A.Murari, J.Vega, G.Vagliasindi, J.A.Alonso, D.Alves, R.Coelho, S.DeFiore, J.Farthing, C.Hidalgo, G.A.Rattá, and JET-EFDA Contributors. “Recent Developments in Data Mining and Soft Computing for JET with a view on ITER”. Fusion Engineering and Design Volume 84, Issues 7-11, Pages 1372-1375. Junio 2009.
- [72] A. Santagiustina, S.A: Arshad, G. Bosia, M. Browne, D.J. Campbell, Dapos, G. Antona, G. Fishpool, G.F Neill. “Design of the m=2, n=1 tearing mode control system for JET”. Fusion Engineering - Supplement, 1993., 15th IEEE/NPSS Symposium on Volume, Issue , 11-15. Page(s):58 – 61. 1993.
- [73] B.Cannas, A.Fanni, G.Sias, P.Sonato, M.K. Zedda. “Neural approaches to disruption prediction at JET”. 31st EPS Conference on Plasma Phys. London, ECA Vol.28G, P-1.167. 28 Junio - 2 Julio. 2004.
- [74] B. Cannas, et al. “A prediction tool for real-time application in the disruption protection system at JET”. Nucl. Fusion 47 No 11. 1559-1569. Noviembre 2007.
- [75] G. A. Rattá, J. Vega, A. Murari, G. Vagliasindi, M. F. Johnson, P.C. de Vries and JET EFDA Contributors. “An Advanced Disruption Predictor for JET tested in a

simulated Real Time Environment” Aceptado para su publicación en Nucl. Fusion. 2009.

[76] A.Murari, J.Vega, D.Mazon, G.A. Ratta, J.Svensson and G.Vagliasindi C. Boulbe, B. Faugeras and JET-EFDA Contributors. “New Information Processing Methods for Control in Fusion.” (To be submitted to Fusion Eng. & Design). 2009.

[77] A.Murari, J.Vega, D.Mazon, G. A. Rattá, J.Svensson, G.Vagliasindi, D.Alves, P.Arena, C. Boulbe, R.Coelho, B. Faugeras, L.Fortuna, D. Moreau, D.Testa and JET-EFDA Contributors. “Innovative Signal Processing and Data Analysis Methods for Control in reactor relevant devices”. Será enviado a Nuclear Fusion. 2009.

[78] P.C. de Vries, M.F. Jonson and I. Segui. “Statistical Analysis of Disruptions in JET”. Nucl. Fusion 49 055011. 12pp. 2009.

[79] T. Briggs. 2005. MATLAB/MEX Interface to SVMlight. www document, <http://www.ship.edu/~thb/mexsvm/>.

[80] I. Barlow, M. Bigi, J. Bird, G. Bonizzoni, R. Buttery and R. Clay et al., “The error field correction coils on the JET machine”, Fusion Eng. Des. 58–59. pp. 189–193. 2001.

[81] J. Vega, A. Murari, G. Rattá, S. González, S. Dormido-Canto and JET-EFDA Contributors. “Progress on statistical learning systems as data mining tools for the creation of automatic databases in Fusion environments”. (Será enviado a Fusion Eng. & Design) 2009.

[82] J. Vega, A. Murari, S. González y JET-EFDA Contributors. “An universal method for automatic event location in waveforms and video-movies: applications to massive nuclear fusion databases”. Enviado a Review of Scientific Instruments. 2009.

[83] V. Vovk, A. Gammerman y G. Shafer. “Algorithmic learning in a random world”. New York: Springer, 2005.



---

## Apéndice

---

La presente Memoria contiene, parcialmente, el trabajo presentados en los siguientes artículos y comunicaciones.

### Publicaciones en revistas

J. Vega, G. A. Rattá, A. Murari, P. Castro, S. Dormido-Canto, R. Dormido, G. Farias, A. Pereira, A. Portas, E. De la Luna, I. Pastor, J. Sánchez, N. Duro, R. Castro, M. Santos, H. Vargas. *“Recent results on structural pattern recognition for Fusion massive databases”*. Proc. of the IEEE International Symposium on Intelligent Signal Processing. ISBN: 1-4244-0829-6 (2007) 949-954.

A. Murari, J. Vega, J. A. Alonso, E. de la Luna, J. Farthing, C. Hidalgo, G. A. Rattá, J. Svensson, G. Vagliasindi and JET EFDA Contributors. *“New techniques and technologies for information retrieval and knowledge extraction from nuclear fusion massive databases”*. Proc. of the IEEE International Symposium on Intelligent Signal Processing. ISBN: 1-4244-0829-6 (2007) 943-948.

- J. Vega, A. Murari, G. A. Rattá, P. Castro, A. Pereira, A. Portas and JET-EFDA Contributors. *“Structural pattern recognition techniques for data retrieval in massive fusion databases”*. Burning Plasma Diagnostics. AIP Conference Proceedings (ISBN 978-0-7354-0507-3). 988 (2008) 481-484.
- A. Murari, J. Vega, J. A. Alonso, E. de la Luna, J. Farthing, C. Hidalgo, G. A. Rattá, J. Svensson, G. Vagliasindi and JET EFDA Contributors. *“How to extract information and knowledge from fusion massive databases”*. Burning Plasma Diagnostics. AIP Conference Proceedings (ISBN 978-0-7354-0507-3). 988 (2008) 457-470.
- G. A. Rattá, J. Vega, A. Pereira, A. Portas, E. De la Luna, S. Dormido-Canto, G. Farias, R. Dormido, J. Sánchez, N. Duro, H. Vargas, M. Santos, G. Pajares, A. Murari and JET EFDA Contributors. *“First applications of structural pattern recognition methods to the investigation of specific physical phenomena at JET”*. Fusion Engineering and Design. 83 (2008) 467-470.
- S. Dormido-Canto, G. Farias, R. Dormido, J. Vega, J. Sánchez, N. Duro, H. Vargas, G. Rattá, A. Pereira, A. Portas. *“Structural pattern recognition methods based on string comparison for fusion databases”*. Fusion Engineering and Design. 83 (2008) 421-424.
- J. Vega, G. A. Rattá, P. castro, A. Murari and JET-EFDA Contributors. *“Development of learning systems with data tours techniques for fusion databases”*. Computational Intelligence in Decision and Control. Proc.of the 8th International FLINS Conference. 1 (2008) 103-108.
- G. de Arcas, J. M. López, M. Ruiz, E. Barrera, J. Nieto, J. Veja, G. A. Rattá, A. Murari. *“Design of an advanced intelligent instrument with waveform recognition based on the ITMS platform”*. Computational Intelligence in Decision and Control. Proc.of the 8th International FLINS Conference. 1 (2008) 423-428.
- G. A. Rattá, J. Vega, A. Murari, M. Johnson and JET-EFDA Contributors. *“Feature extraction for improved disruption prediction analysis at JET”*. Review of Scientific Instruments. 79, 10F328 (2008).

J. Vega, A. Murari, A. Pereira, A. Portas, G. A. Rattá, R. Castro and JET-EFDA Contributors. *“Overview of intelligent data retrieval methods for waveforms and images in massive fusion databases”*. Fusion Engineering and Design.  
doi:10.1016/j.fusengdes.2008.11.097

A.Murari, J.Vega, G.Vagliasindi, J.A.Alonso, D.Alves, R.Coelho, S.DeFiore, J.Farthing, C.Hidalgo, G.A.Rattá, and JET-EFDA Contributors. *“Recent Developments in Data Mining and Soft Computing for JET with a view on ITER”*. Fusion Engineering and Design. Volume 84, Issues 7-11, Junio 2009, Pages 1372-1375  
doi:10.1016/j.fusengdes.2008.12.060

G. A. Rattá, J. Vega, A. Murari, G. Vagliasindi and M. Johnson and JET-EFDA Contributors. *“An Advanced Disruption Predictor for JET tested in a simulated Real Time Environment”*. Aceptado para su publicación en Nuclear Fusion.

J.Vega, A. Murari, G. Vagliasindi, G. A. Rattá and JET-EFDA Contributors. *“Automated estimation of L/H transition times at JET by combining Bayesian statistics and Support Vector Machines”*. 2009. Nucl. Fusion 49 085023 (11pp)  
doi:10.1088/0029-5515/49/8/085023

A.Murari, J.Vega, G.A.Rattá, G.Vagliasindi, M.F.Johnsons, S.H.Hong and JET-EFDA Contributors. *“Unbiased and Non-Supervised Learning Methods for Disruption Prediction at JET”*. Nucl. Fusion 49 (2009) 055028 (11pp) .

M. Ruiz, J. Vega, E. Barrera, J. González, A. Murari, R. Meléndez, G.Rattá, S. González and JET-EFDA Contributors. *“Test-bed of a real time detection system of L-H & H-L plasma transitions implemented with ITMS platform”*. Enviado a Fusion Eng. & Design.

J. Vega, A. Murari, G. Rattá, S. González, S. Dormido-Canto and JET-EFDA Contributors. *“Progress on statistical learning systems as data mining tools for the creation of automatic databases in Fusion environments”*. Enviado a Fusion Eng. & Design.

A.Murari, J.Vega, D.Mazon, G.A. Ratta, J.Svensson, G.Vagliasindi C. Boulbe, B. Faugeras. *“New Information Processing Methods for Control in Fusion”*. Enviado a Fusion Eng. & Design.

A.Murari, J.Vega, D.Mazon, G. A. Rattá, J.Svensson, G.Vagliasindi, D.Alves, P.Arena, C. Boulbe, R.Coelho, B. Faugeras, L.Fortuna, D. Moreau, D.Testa and JET-EFDA Contributors. *“Innovative Signal Processing and Data Analysis Methods for Control in reactor relevant devices”*. Enviado a Nuclear Fusion.

A. Murari, J.Vega, D. Patané, G.Vagliasindi, D.Mazon, N.Martin, N.F.Martin, G.Rattá, P.Arena, V.Caloone5 and JET-EFDA Contributors. *“Machine Learning for the Identification of Scaling Laws and Dynamical Systems directly from data in Magnetic Confinement Fusion”*. Será enviado a Nuclear Instruments and Methods (the FDT1 conference, Frascati, Italia).

J. Sánchez,...., G.Rattá,...., and B. Zurro. *“Confinement transitions in TJ-II under Li-coated wall conditions”*. Nuclear Fusion 49, 104018 (2009).

Iop:stacks.iop.org/NF/49/104018

### Comunicaciones a congresos

J. Vega, R. Castro, A. Portas, A. Pereira, P. Castro, G. Rattá, J. Sánchez, A. Pastoriza, P. Manrubia, F. Sastre, M. Ruiz, E. Barrera. *“Los nuevos paradigmas de interacción con los entornos experimentales de dispositivos de fusión”*. 33ª Reunión Anual de la Sociedad Nuclear Española. 26-28 Septiembre 2007.

G. A. Rattá, J. Vega, A. Pereira, A. Portas, E. De la Luna, S. Dormido-Canto, G. Farias, R. Dormido, J. Sánchez, N. Duro, H. Vargas, M. Santos, G. Pajares, A. Murari and JET EFDA Contributors. *“First applications of structural pattern recognition methods to the investigation of specific physical phenomena at JET”*. 6th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research. 4-8 Junio 2007. Inuyama (Japón) (<http://tm2007.nifs.ac.jp/>).



S. Dormido-Canto, G. Farias, R. Dormido, J. Sánchez, N. Duro, H. Vargas, J. Vega, G. Rattá, A. Pereira, A. Portas. *“Comparison of structural pattern recognition methods for fusion databases”*. 6th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research. 4-8 Junio 2007. Inuyama (Japan) (<http://tm2007.nifs.ac.jp/>).

D. Raju, J. Vega, P. Castro, G. A. Rattá, A. Murari, G. Vagliasindi and JET EFDA Contributors. *“Structural pattern recognition for image processing in fusion plasmas”*. 6th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research. 4-8 Junio 2007. Inuyama (Japan) (<http://tm2007.nifs.ac.jp/>).

A. Pastoriza, P. Manrubia, F. Sastre, M. Ruiz, E. Barrera, J. Vega, R. Castro, A. Portas, A. Pereira, P. Castro, G. Rattá. *“On-line reprogramming of data acquisition systems during shots”*. 6th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research. 4-8 Junio 2007. Inuyama (Japan) (<http://tm2007.nifs.ac.jp/>).

J. Vega, A. Murari, G. A. Rattá, P. Castro, A. Pereira, A. Portas and JET EFDA Contributors. *“Structural pattern recognition techniques for data retrieval in massive fusion databases”*. International Workshop on Burning Plasma Diagnostics. Villa Monastero, Varenna, Italia. 24-28 Septiembre 2007.

A. Murari, J. Vega, J. A. Alonso, E. de la Luna, J. Farthing, C. Hidalgo, G. A. Rattá, J. Svensson, G. Vagliasindi and JET EFDA Contributors. *“How to extract information and knowledge from fusion massive databases”*. International Workshop on Burning Plasma Diagnostics. Villa Monastero, Varenna, Italia. 24-28 Septiembre 2007.

J. Vega, G. Rattá, A. Murari, P. Castro, S. Dormido-Canto, R. Dormido, G. Farias, A. Pereira, A. Portas, E. de la Luna, I. Pastor, J. Sánchez, N. Duro, R. Castro, M. Santos, H. Vargas and JET EFDA Contributors. *“Recent results on structural pattern recognition for Fusion massive databases”*. IEEE International Symposium on Intelligent Signal Processing. Alcalá de Henares (Madrid) España. 3-5 Octubre 2007. Proc. of the IEEE International Symposium on Intelligent Signal Processing. ISBN: 1-4244-0829-6 (2007) 949-954.

A. Murari, J. Vega, J. A. Alonso, E. de la Luna, J. Farthing, C. Hidalgo, G. A. Rattá, J. Svensson, G. Vagliasindi and JET EFDA Contributors. *“New techniques and technologies for information retrieval and knowledge extraction from nuclear fusion massive databases”*. IEEE International Symposium on Intelligent Signal Processing. Alcalá de Henares (Madrid) España. 3-5 Octubre 2007. Proc. of the IEEE International Symposium on Intelligent Signal Processing. ISBN: 1-4244-0829-6 (2007) 943-948.

G. A. Rattá, J. Vega, A. Murari, M. Johnson and JET-EFDA Contributors. *“Feature extraction for improved disruption prediction analysis at JET”*. 17th Topical Conference on High Temperature Plasma Diagnostics. May 11-15, 2008. Albuquerque. New Mexico. (USA). <http://www.esc.sandia.gov/httpd08/httpd08ConferenceInfo.html>

J. Vega, A. Murari, A. Pereira, A. Portas, G. A. Rattá, R. Castro and JET-EFDA Contributors. *“Overview of intelligent data retrieval methods for waveforms and images in massive fusion databases”*. 25th Symposium on Fusion Technology. 15th – 19th Septiembre 2008. Rostock (Alemania).  
<http://www.ipp.mpg.de/eng/for/veranstaltungen/soft2008/>

A. Murari, J. Vega, G. Vagliasindi, J.A. Alonso, D. Alves, R. Coelho, S. DeFiore, J. Farthing, C. Hidalgo, G.A. Rattá, and JET-EFDA Contributors. *“Recent Developments in Data Mining and Soft Computing for JET with a view on ITER”*. 25th Symposium on Fusion Technology. 15th – 19th Septiembre 2008. Rostock (Alemania).  
<http://www.ipp.mpg.de/eng/for/veranstaltungen/soft2008/>

J. Vega, G. A. Rattá, P. Castro, A. Murari and JET-EFDA Contributors. *“Development of learning systems with data tours techniques for fusion databases”*. 8th International FLINS Conference (<http://www.mat.ucm.es/congresos/flins2008/>). 21st – 24th Septiembre 2008. Madrid (España).

G. de Arcas, J. M. López, M. Ruiz, E. Barrera, J. Nieto, J. Veja, G. A. Rattá, A. Murari. *“Design of an advanced intelligent instrument with waveform recognition based on the ITMS platform”*. 8th International FLINS Conference. 21 – 24 Septiembre 2008. Madrid (España). <http://www.mat.ucm.es/congresos/flins2008/>

M. Ruiz, J. Vega, E. Barrera, J. González, A. Murari, R. Meléndez, G. Rattá, S. González and JET-EFDA Contributors. *“Test-bed of a real time detection system of L-H & H-L plasma transitions implemented with ITMS platform”*. 7th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research. 15-19 Junio 2009. Aix-en-Provence. Francia.

J. Vega, A. Murari, G. Rattá, S. González, S. Dormido-Canto and JET-EFDA Contributors. *“Progress on statistical learning systems as data mining tools for the creation of automatic databases in Fusion environments”*. 7th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research. 15-19 Junio 2009. Aix-en-Provence. Francia.

G. A. Rattá, J. Vega, A. Murari, G. Vagliasindi. *“Inspection of Disruptive Behaviours at JET using Generative Topographic Maps”*. 4th International Scientific Conference on Physics and Control (PHYSCON 2009). Catania, Italia.

A. Murari, J. Vega, D. Mazon, G. A. Ratta, J. Svensson, G. Vagliasindi, C. Boulbe, B. Faugas. *“New Information Processing Methods for Control in Fusion”*. 7th IAEA Technical Meeting on Control, Data Acquisition and Remote Participation for Fusion Research. 15-19 Junio 2009. Aix-en-Provence. Francia.

### **Seminarios**

G. A. Rattá, J. Vega, A. Murari, G. Vagliasindi and M. Johnson and JET-EFDA Contributors. *“An Advanced Disruption Predictor for JET tested in a simulated Real Time Environment”*. JET Task Force D Meeting, Culham science Center 11 de diciembre de 2008.