



UNIVERSITAT OBERTA DE CATALUNYA (UOC)
MÁSTER UNIVERSITARIO EN CIENCIA DE DATOS (*Data Science*)

TRABAJO FINAL DE MÁSTER

ÁREA: 2

Detección y clasificación de eventos anómalos en la extracción de petróleo

Métodos de aprendizaje incremental aplicados al dataset 3W de Petrobras

Autor: Baltasar Boix Rita

Tutor: Jesús López Lobo

Profesor: Jordi Casas Roma

Tarragona, 4 de enero de 2023

Créditos/Copyright

Dataset 3W de Petrobras [34] in github licensed under the Apache License 2.0:

”This is the first Petrobras’ repository on GitHub. It supports the 3W project and promotes experimentation of Machine Learning-based approaches and algorithms for specific problems related to undesirable events that occur in offshore oil wells.”



Esta obra está sujeta a una licencia de Reconocimiento - NoComercial - SinObraDerivada
3.0 España de Creative Commons.

FICHA DEL TRABAJO FINAL

Título del trabajo:	Detección y clasificación de eventos anormales en la extracción de petróleo.
Nombre del autor:	Baltasar Boix Rita
Nombre del colaborador docente:	Jesús López Lobo
Nombre del PRA:	Jordi Casa Roma
Fecha de entrega (mm/aaaa):	01/2023
Titulación o programa:	Máster Universitario en Ciencia de Datos
Área del Trabajo Final:	TFM Área 2
Idioma del trabajo:	Español
Palabras clave	Fault detection and diagnosis, Oil well monitoring, Abnormal event management, Multivariate time series classification, Stream learning.

Cita

Learn from the mistakes of others. You can't live long enough to make them all yourself.

— Eleanor Roosevelt

Agradecimientos

A mi familia.

Abstract/Resumen

La temprana detección de eventos indeseables en pozos de petróleo y gas puede ayudar a prevenir accidentes con afectación a las personas y al medio ambiente además de reducir pérdidas de producción y costes de mantenimiento. Se ha utilizado el dataset 3W de Petrobras que contiene series temporales etiquetadas por expertos en condiciones normales de operación y en 8 tipos de situaciones indeseables. Se ha emulado la recepción de datos en tiempo real y se han tratado con modelos de *stream learning*. Se ha demostrado que es posible llevar a cabo una predicción en tiempo real de los eventos anómalos. Asimismo, se ha evaluado la presencia del fenómeno concept drift durante el proceso de aprendizaje con diferentes métodos a partir de los datos reales del dataset.

The early detection of undesirable events in oil and gas wells can help to prevent accidents with affectation to people and the environment in addition to reducing production losses and maintenance costs. The 3W dataset of Petrobras have been used. It contains time series labelled by experts in normal and in 8 types of undesirable situations. The reception of data in real time has been emulated and treated with *stream learning* models. It has been demonstrated that it is possible to predict anomalous events in real-time. Concept drift phenomenon has been evaluated during the learning process with different methods applied to the dataset real data.

Palabras clave: Fault detection and diagnosis, Oil well monitoring, Abnormal event management, Multivariate time series classification, stream learning.

Índice general

Abstract/Resumen	IX
Índice	XI
Llistado de Figuras	XIII
Listado de Tablas	1
1. Introducción	3
1.1. Descripción general del problema	3
1.2. Motivación personal	4
1.3. Objetivos del TFM	5
1.4. Descripción de la metodología	6
1.5. Planificación	6
2. Estado del arte	9
2.1. Online Learning	9
2.2. Concept Drift	10
2.3. Outliers	12
2.4. Trabajos sobre <i>3W Dataset</i>	13
3. Desarrollo del TFM	15
3.1. Descripción del Dataset 3W	15
3.2. Análisis exploratorio de las series reales	16
3.3. Preprocesamiento de los datos	20
3.4. Métodos Offline	20
3.4.1. Redes Neuronales	20
3.4.2. Modelos de Sklearn de aprendizaje incremental	21
3.5. Métodos Online	23
3.5.1. Modelo base	23

3.5.2. Detección de <i>drift</i>	26
4. Conclusiones del TFM	35
Bibliografía	35

Índice de figuras

1.1. Planning del Trabajo Final de Master	7
2.1. Batch Learning[6]	10
2.2. Online Learning[6]	10
2.3. Orígenes de <i>Concept Drift</i> [23]	11
2.4. Tipos de <i>Concept Drift</i> [23]	11
2.5. Esquema de un detector de <i>Concept Drift</i> [23]	11
2.6. Esquema genérico para un <i>online adaptive learning algorithm</i> . [15]	11
2.7. Taxonomía de métodos de detección no supervisada de <i>concept drift</i> [16]	12
2.8. Diagrama de bloques de la aplicación[26]	13
2.9. Metodología propuesta por Aslam <i>et al</i> [2]	14
3.1. Esquema simplificado de un típico pozo marino de flujo natural.[34]	16
3.2. Ejemplo de series temporales del dataset 3W	17
3.3. Distribución de las series temporales del dataset 3W	18
3.4. Cronograma de la adquisición de datos reales del dataset 3W	19
3.5. Redes neuronales ensayadas	21
3.6. Evolución del aprendizaje de la RNN	22
3.7. Resultados de la RNN para el conjunto de test	22
3.8. Resultado de aplicar ipca+RF al conjunto de test	23
3.9. Diagrama del modelo online utilizado	24
3.10. Modelo Hoeffding Tree para la detección de <i>Flow Instability</i> con datos reales	25
3.11. HT Train Accuracy	25
3.12. HT Evaluation Accuracy	25
3.13. Precisión del Hoeffding Tree	25
3.14. ALMA Train Accuracy	26
3.15. ALMA Evaluation Accuracy	26
3.16. Precisión del linear_model.ALMAClassifier	26

3.17. Comparación del modelo Hoeffding Tree con el NoChange para <i>Flow Instability</i> con datos reales	27
3.18. Comparación del modelo Hoeffding Tree con el NoChange para <i>Rapid Productivity Loss</i> con datos reales y simulados ordenados aleatoriamente	28
3.19. Diagrama del modelo online utilizado con detección de <i>drift</i>	28
3.20. Detección de <i>drift</i> con diferentes métodos para <i>flow instability</i>	32
3.21. Detección de <i>drift</i> con ADWIN para anomalías 1 (ABRUPT INCREASE OF BSW), 2 (SPURIOUS CLOSURE OF DHSV) y 3 (SEVERE SLUGGING) . . .	33
3.22. Detección de <i>drift</i> con ADWIN para la anomalía 5 (<i>Rapid Productivity Loss</i>) con series simuladas	34

Índice de cuadros

3.1. <i>Tamaño estimado de las ventanas de tiempo utilizadas para confirmar la ocurrencia de eventos no deseados</i> [34]	18
3.2. Porcentaje de Nan en el Dataset 3W	18

Capítulo 1

Introducción

1.1. Descripción general del problema

El uso de ordenadores para el control de procesos complejos ha permitido grandes avances. Acciones de bajo nivel como la apertura y cerrado de válvulas, conocido como control regulatorio, que anteriormente era realizado por operadores humanos son ahora actuados automáticamente. Aun así, la respuesta ante situaciones anormales del proceso sigue estando en buena medida en manos de los operadores. Esto implica la detección temprana de las situaciones anormales, diagnosticando sus causas y tomando las acciones de control apropiadas para llevar el proceso de nuevo a un estado normal y seguro. Al conjunto de esta actividad se le ha denominado *Abnormal Event Management* (AEM), un componente clave del control supervisor[35]. Las estadísticas de Seguridad de la Industria de Proceso muestran que los grandes accidentes son infrecuentes pero cuando ocurren pueden transformar a todo un sector (*BP Deepwater Horizon Oil Spill* [28]). Los accidentes menores son más comunes pero también causan lesiones y enfermedades laborales o afectaciones medioambientales además de un gran coste acumulado. El primer paso en AEM es la automatización de la detección y diagnóstico de fallos en el proceso industrial. Para resolver este problema se han utilizado diferentes métodos que Venkatasubramanian et al.(2003)[35] clasifican principalmente como los basados en modelos y los basados únicamente en datos históricos del proceso industrial.

Se utilizan técnicas de aprendizaje automático (*Machine Learning*, ML) para entrenar modelos que detectan y diagnostican situaciones anómalas a partir de datos en tiempo real. Podemos distinguir dos formas de entrenar los modelos[6]:

- Batch/Offline learning. Se selecciona un dataset de datos y se aplican técnicas de aprendizaje supervisado. Periódicamente se actualiza el dataset y se reentrena el modelo.
- Stream/Online learning. Los datos se adquieren secuencialmente en línea y se utilizan para

actualizar el modelo. Los métodos de ML utilizados permiten el aprendizaje incremental. Es ideal para los sistemas que reciben datos en tiempo real y necesitan poder adaptarse a condiciones que cambian rápidamente *Concept Drift*[32].

El efecto *Concept Drift* consiste en el cambio que se puede producir en la distribución de los datos. Estos cambios pueden ser de varios tipos: súbitos, graduales, incrementales o cíclicos[23]. Si esto ocurre a menudo los modelos entrenados quedan obsoletos muy rápidamente. Se hace necesario detectar estos cambios en tiempo real y reentrenar los modelos.

En este TFM utilizaremos métodos de *stream learning*[27] aplicados a las series temporales del dataset 3W de petrobras[34].

El uso de *stream learning* en este caso presenta varios retos:

- El etiquetado lo realizan expertos en batch observando ventanas de tiempo de 5 minutos a 12 horas según la situación analizada (Tabla 3.1). O sea, la observación de un solo registro no contiene suficiente información para clasificarlo y será necesario analizar tendencias.
- El origen de los datos de cada serie temporal es diverso. Hay que pretratarlos para poder utilizar el conjunto.
- Los creadores del dataset[34] advierten que algunos de los tipos de situaciones anómalas pueden sufrir de *concept drift*, cambios en el comportamiento de los datos que ocurren con el tiempo y deterioran el rendimiento de los modelos ajustados (Krawczyk et al., 2017[21]).

1.2. Motivación personal

En mis más de 30 años de experiencia en actividades de Producción de una gran instalación industrial me he enfrentado demasiadas veces a los problemas derivados de situaciones operatorias anómalas (averías de grandes máquinas, cortes de suministro eléctrico, roturas de calderas...). Expongo un caso real que viví en primera persona:

En la sala de reuniones de Dirección estaban los responsables del negocio que tenían algo que decir sobre el incidente de la semana anterior en el que varios hornos de etileno habían sufrido temperaturas que provocaron que se fundieran los serpentines; el coste de reparación y el lucro cesante ascendía a varios millones de euros;

afortunadamente no hubo que lamentar ningún accidentado. Entre todos los presentes acreditábamos más de 100 años de experiencia en la operación de Unidades de Proceso. El jefe de Planta expone la línea temporal de eventos y la evolución de las variables de Proceso más significativas. Se discute las decisiones tomadas por el operador de Control durante la emergencia. Se genera una discusión de pros y contras de las distintas opciones que nos lleva unas dos horas: se podría haber hecho mejor y se habría evitado buena parte del coste!

El operador de Control durante la emergencia estaba solo en su consola en una sala de control donde las alarmas no paraban de sonar y tuvo dos minutos para decidir.

Disponer de un programa de diagnóstico que ayude a saber que está pasando durante una emergencia en tiempo real nos puede ayudar a evitar las peores consecuencias de un incidente y añadimos una capa más al sistema de Seguridad[30]. El operador de Control también es una capa de Seguridad: ni la única ni la más importante.

1.3. Objetivos del TFM

En este trabajo se pretende demostrar que es posible llevar a cabo una predicción en tiempo real de los eventos anómalos producidos en la extracción de crudo en el dataset 3W de Petrobras [34].

- Objetivo principal. Contraste de hipótesis:
 - Hipótesis nula (H_0): no se puede realizar una predicción en tiempo real fiable de los eventos anómalos producidos en la extracción de crudo.
 - Hipótesis alternativa (H_1): es posible realizar una predicción en tiempo real fiable de los eventos anómalos producidos en la extracción de crudo
- Objetivos secundarios:
 1. Considerar de forma adecuada la presencia de outliers durante el proceso de aprendizaje.
 2. Considerar de forma adecuada la presencia del fenómeno concept drift durante el proceso de aprendizaje.

1.4. Descripción de la metodología

Se utilizan principalmente métodos de *stream learning* del paquete River[\[27\]](#) recogiendo métricas de rendimiento de modelos en tiempo real .

- Análisis exploratorio y descriptivo del dataset
- Preprocesado de los datos. Reducción del volumen de datos por utilización de medias por minuto. Normalización de las variables para conjugar las series de varios orígenes.
- Aplicar algoritmos de aprendizaje incremental (*base learner*). Compararlos estadísticamente.
- Detección de concept drift con varias técnicas. Introducir drift artificialmente[\[31\]](#).
- Comparar combinaciones detector de drift/*(base learner)*. Realizar estudio estadístico para confirmar la significancia.

1.5. Planificación

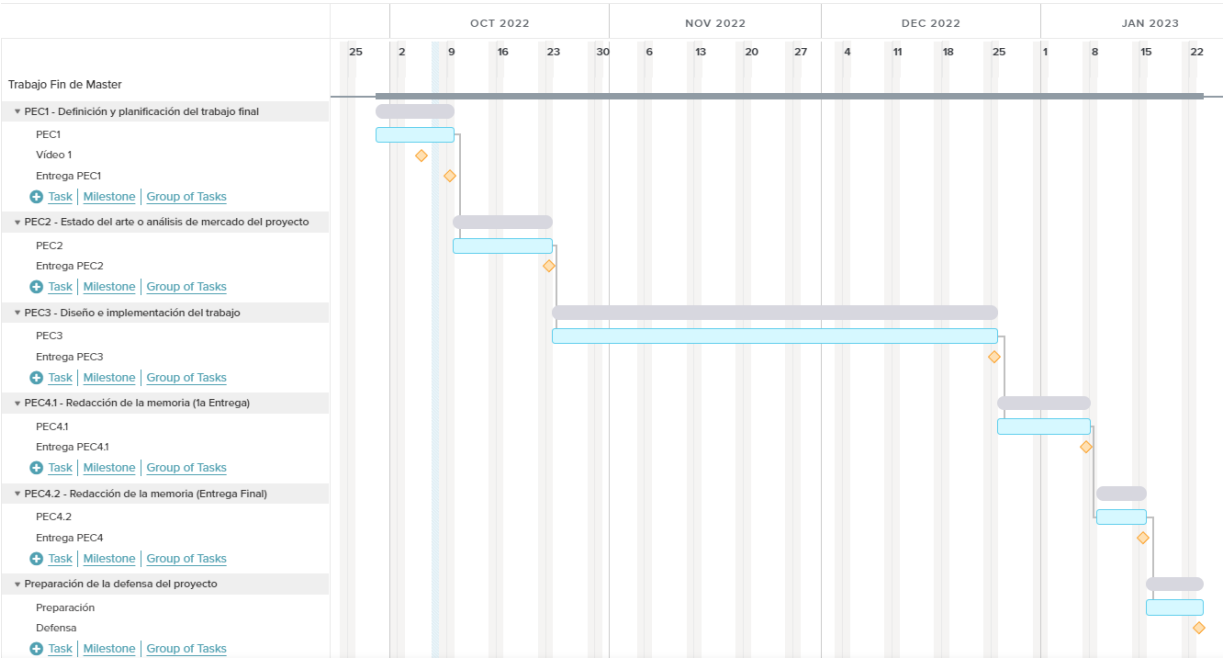


Figura 1.1: Planning del Trabajo Final de Master

Capítulo 2

Estado del arte

2.1. Online Learning

Podemos distinguir dos formas de entrenar los modelos[6]:

- Batch/Offline learning. El modelo de ML se entrena utilizando todo el conjunto de datos que está disponible en un momento determinado. Una vez que tenemos un modelo que funciona bien en el conjunto de prueba, el modelo se envía a producción y, por lo tanto, finaliza el aprendizaje. Si con el paso del tiempo están disponibles nuevos datos es necesario actualizar el modelo. El modelo se entrena desde cero nuevamente utilizando las muestras de datos anteriores y las muestras de datos nuevos.
- Stream/Online learning. Los datos se transmiten (ya sea individualmente o en mini lotes) al algoritmo de ML y actualizan el modelo. El aprendizaje en línea es ideal en situaciones en las que los datos se generan continuamente en el tiempo y necesitamos usar muestras de datos en tiempo real para construir un modelo de predicción. Un ejemplo típico de este caso es la predicción del mercado de valores.

El entorno para *Online Learning* impone unas restricciones computacionales con respecto a la configuración tradicional del *Batch Learning*[22]:

- Cada muestra se procesa solo una vez a su llegada, y los modelos deben poder procesar las muestras secuencialmente tan pronto como se reciben
- El procesamiento de cada muestra debe realizarse dentro del intervalo de tiempo entre muestras
- El algoritmo debe usar solo una cantidad de memoria preasignada y finita

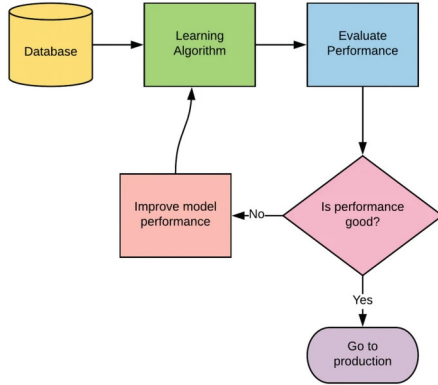


Figura 2.1: Batch Learning[6]

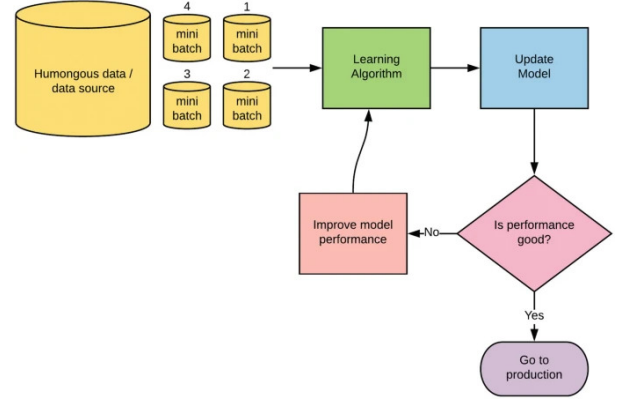


Figura 2.2: Online Learning[6]

- Tras cada exploración de datos, el modelo resultante debe ser válido para sus uso
- El modelo producido debe ser equivalente al que se construiría por *Batch Learning*.

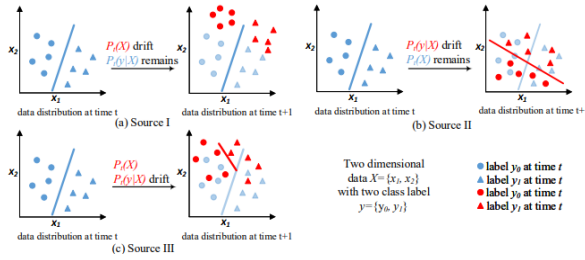
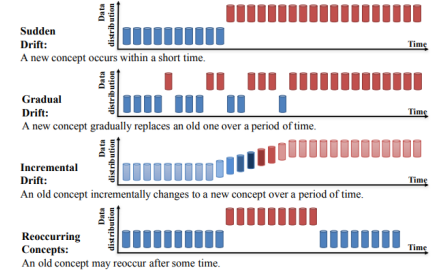
Gomes *et al.*[17] hacen un repaso de los retos a los que se enfrenta el *Stream Learning* y las ventajas que tiene al enfrentarse al Big Data caracterizado por las 3 V's(Volumen, Variedad, Velocidad).

En este TFM utilizaremos métodos de *stream learning* del módulo de Python River[27].

2.2. Concept Drift

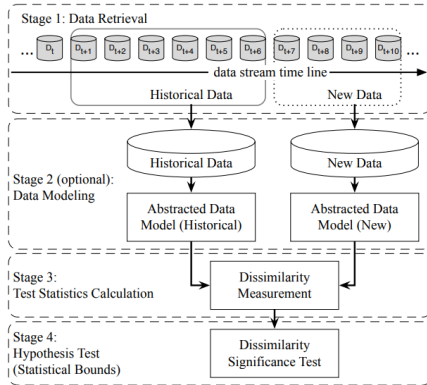
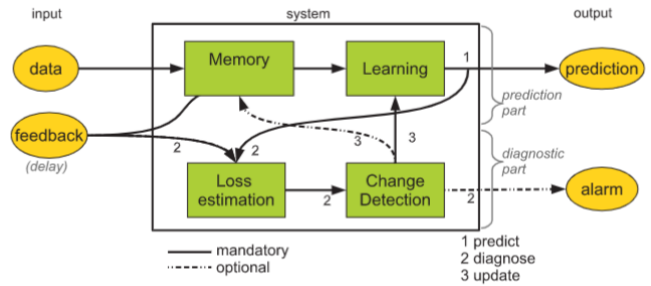
Lu *et al.*[23] consideran que *Concept drift* en el tiempo t puede definirse como el cambio en la probabilidad conjunta de X e y en el tiempo t . Puesto que la probabilidad conjunta $P_t(X, y)$ puede descomponerse en dos partes: $P_t(X, y) = P_t(X)P_t(y|X)$, *concept drift* puede provenir de tres orígenes (Figura 2.3:

- Origen I: $P_t(X) \neq P_{t+1}(X)$ mientras que $P_t(y|X) = P_{t+1}(y|X)$, esto es, varía $P_t(X)$ mientras que $P_t(y|X)$ no cambia. Ya que $P_t(X)$ no afecta la frontera de decisión, se considera *virtual drift*.
- Origen II: $P_t(y|X) \neq P_{t+1}(y|X)$ mientras que $P_t(X) = P_{t+1}(X)$. En este caso la frontera de decisión cambia y por tanto la precisión de los modelos empeora, se considera *Real Drift*.
- Origen III: mezcla del Origen I y II, o sea, $P_t(X) \neq P_{t+1}(X)$ y $P_t(y|X) \neq P_{t+1}(y|X)$.

Figura 2.3: Orígenes de *Concept Drift*[23]Figura 2.4: Tipos de *Concept Drift*[23]

Se pueden distinguir varios tipos de *Concept Drift*: súbitos, graduales, incrementales o cíclicos[23] (Figura 2.4).

Uno de los principales retos de los métodos para detectar *Concept Drift* es distinguirlo del ruido en el flujo de datos[32]. Lu *et al.*[23] proponen un marco general (Figura 2.5) para un detector de *Concept Drift* y hacen un repaso de diferentes algoritmos. Gama *et al.*[15] tratan el problema de adaptar los métodos de *Stream Learning* a la presencia de *Concept Drift* y proponen el esquema general de la Figura 2.6.

Figura 2.5: Esquema de un detector de *Concept Drift*[23]Figura 2.6: Esquema genérico para un *online adaptive learning algorithm*. [15]

Sethi *et al.*[31] exponen que la detección de *Concept Drift* se ha abordado tradicionalmente como una tarea supervisada, con datos etiquetados que se utilizan constantemente para validar el modelo aprendido. Aunque efectivas, estas técnicas no siempre son prácticas, ya que el etiquetado es una actividad difícil, costosa y que requiere mucho tiempo. Por otro lado, las técnicas de detección de cambios no supervisadas son poco fiables, ya que producen una gran cantidad de falsas alarmas. La ineficacia de las técnicas no supervisadas radica en la exclusión de las características del clasificador del proceso de detección. Hacen un repaso de las técnicas existentes y proponen un método no supervisado.

Hu *et al.*[19] y Gemaque *et al.*[16] hacen un repaso de la técnicas de detección no supervisadas de *Concept Drift*. Gemaque *et al.*[16] proponen la taxonomía de la figura 2.7.

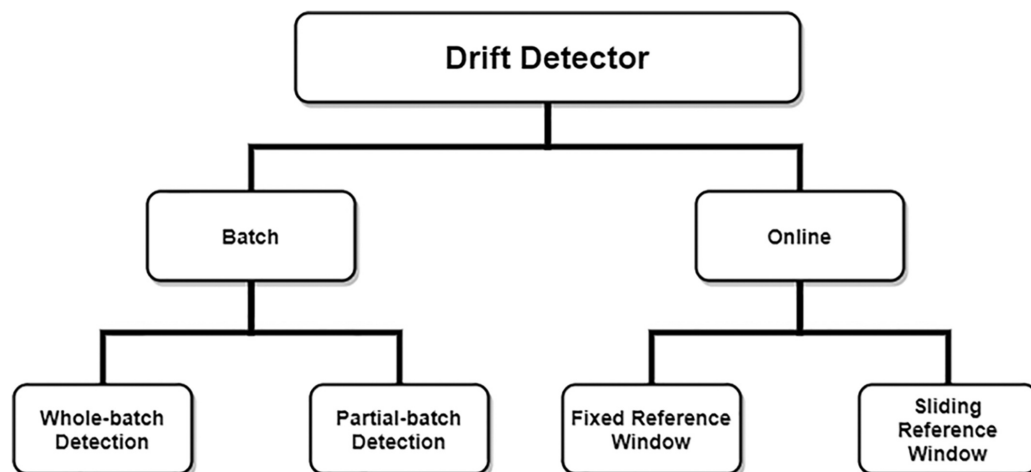


Figura 2.7: Taxonomía de métodos de detección no supervisada de *concept drift*[16]

Lobo *et al.*[24][22] investigan el uso de *Spiking Neural Networks* (SNN) en entornos de *online learning* y su comportamiento frente al *concept drift*.

2.3. Outliers

Gupta *et al.*[18] hacen un extensivo repaso de técnicas de detección de *outliers* clasificándolas en función del origen y tipo de datos e identificando los retos que presenta el tratamiento en tiempo real de un gran flujo de datos. Duraj *et al.*[11] comparan diferentes métodos de detección aplicados a series temporales que explotan la autocorrelación que muestran este tipo de datos. El libro *Outlier Analysis*[1] de Aggarwal presenta en detalle multitud de algoritmos. Pokrajak *et al.*[29] presentan un algoritmo incremental de detección de *outliers* para su uso online.

Hay que resaltar que los algoritmos de detección de *drift* y *outliers* utilizan técnicas muy similares.

2.4. Trabajos sobre *3W Dataset*

Desde la publicación de dataset en 2019[34] se han publicado diversos artículos que han aplicado distintas técnicas de ML para la clasificación de las anomalías. Marins *et al*[26] utilizan un esquema de cálculo (Figura 2.8) que comparten casi todos los artículos. Cada observación se forma a partir de un ventana de tiempo deslizante (*rolling time window*) de la que se calculan una colección de estadísticos (media, desviación estándar, máximo, etc.). Se aplica una reducción del número de variables (*features*) mediante técnicas como PCA (*Principal Component Analysis*). El dataset ya pretratado se alimenta a diferentes algoritmos de clasificación. Los mejores resultados se obtienen utilizando Random Forest (RF) consiguiendo precisiones por encima del 90 %.



Figura 2.8: Diagrama de bloques de la aplicación[26]

Brønstad *et al*[7] refinan la metodología anterior utilizando también RF. Turan y Jäschke[33] proponen el uso de un árbol de decisión (*Decision Tree*) que, aunque obtiene resultados ligeramente inferiores a un RF, es mucho más simple e interpretable. Figueirêdo *et al*[12] tienen un planteamiento diferente, ensayan técnicas de aprendizaje no supervisado (*unsupervised Learning*) para detectar/seleccionar situaciones anómalas que deberán etiquetar los expertos. Carvalho *et al*[9][8] se centran en una de las anomalías, *Flow Instability*, en la selección de variables y en la evaluación de la precisión; alertan del sesgo de similaridad (*similarity bias*) en la selección de los conjuntos de test y entrenamiento. De Salvo *et al*[10] proponen el uso de un gráfico de control de las principales variables (*Control Chart*) para detectar anomalías y utilizan un RF para clasificarlas. Aslam *et al*[2] proponen la metodología de la figura 2.9; la novedad se encuentra en el uso de técnicas de *Explainable Artificial Intelligence* (XAI), que permiten interpretar los modelos caja negra (*black box*). Machado *et al*[25] usan un autoencoder *Long Short-Term Memory* (LSTM) y un *one-class Support Vector Machine* (SVM) para clasificar dos de las anomalías (*Spurious DHSV closure* y *Hydrate in Production Line*): la primera con una dinámica rápida (pocos minutos) y la otra lenta (más de una hora); mejoran la precisión y el F1 de estudios anteriores.

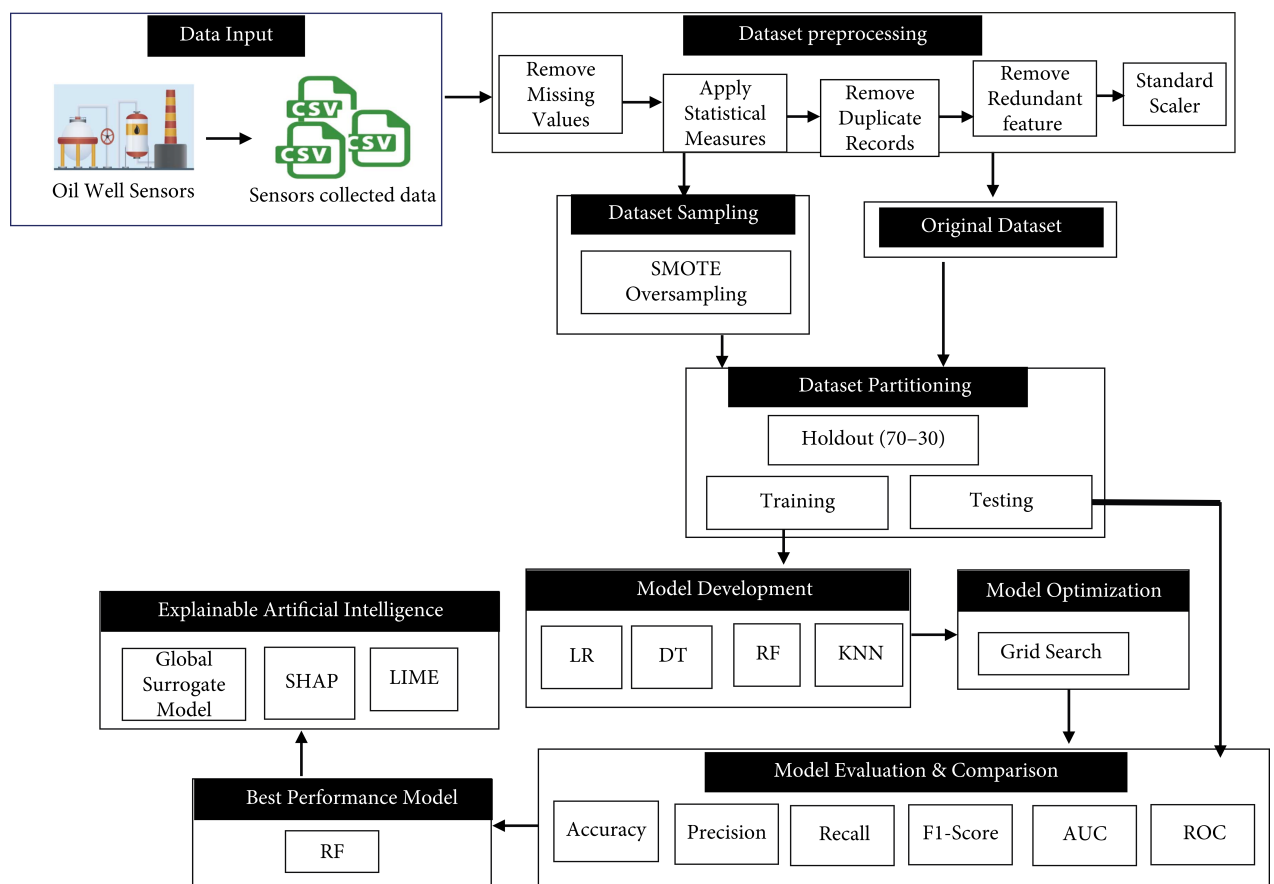


Figura 2.9: Metodología propuesta por Aslam *et al*[2]

Capítulo 3

Desarrollo del TFM

3.1. Descripción del Dataset 3W

En la Figura 3.1 puede verse un esquema simplificado de una instalación típica de extracción de crudo marina.

El dataset 3W[34] está formado por un conjunto de series temporales que describen la operación de un pozo de extracción. Las series tienen como origen:

- Datos reales de pozos. Identifican el pozo con un número.
- Datos generados por simulación.
- Datos manuales generadas por expertos.

Cada serie temporal se almacena en un archivo CSV. Cada observación es una línea en el archivo CSV. La primera línea de cada archivo CSV contiene un encabezado con identificadores de columna. Cada columna de archivos CSV almacena el siguiente tipo de información:

- **timestamp**: Marca de tiempo que actúa como índice
- **P-PDG**: variable de presión en el Medidor Permanente de Fondo de Pozo (PDG) Pa
- **P-TPT**: variable de presión en el Transductor de Temperatura y Presión (TPT) Pa
- **T-TPT**: variable de temperatura en el Transductor de Temperatura y Presión (TPT) °C
- **P-MON-CKP**: variable de presión aguas arriba del regulador de producción (CKP) Pa
- **T-JUS-CKP**: temperatura variable aguas abajo del regulador de producción (CKP) °C
- **P-JUS-CKGL**: variable de presión aguas arriba del regulador de gas-lift (CKGL) Pa

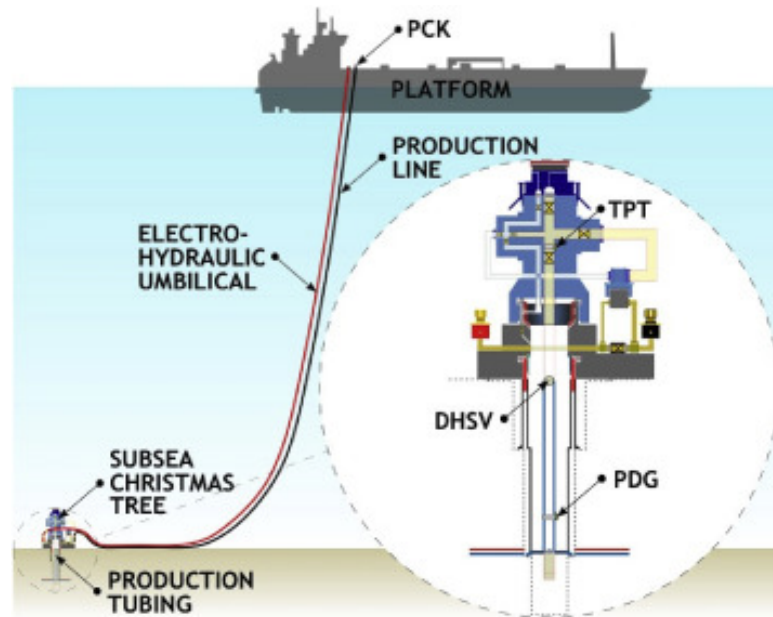


Figura 3.1: Esquema simplificado de un típico pozo marino de flujo natural.[34]

- **T-JUS-CKGL**: temperatura variable aguas arriba del regulador de gas-lift (CKGL) °C
- **QGL**: caudal de gas-lift (inyección de gas para incrementar la producción) m^3/s
- **class**: etiquetas de observaciones asociadas con tres tipos de períodos (normal, transitorio de la anomalía y estado estable defectuoso)

Cada serie está etiquetada a dos niveles por expertos:

- globalmente por la situación normal o anómala que se desarrolla en la serie.
- Por registro (uno por segundo): normal, transición, anómalo.

En el cuadro 3.1 se muestran las situaciones anormales junto con las ventanas de tiempo usadas por los expertos para etiquetarlas. En el apartado 2.2 del artículo de Vargas *et al.*[34] se encuentra una descripción detallada de cada una de ellas.

3.2. Análisis exploratorio de las series reales

En la figura 3.2 se muestra de forma gráfica dos series correspondientes al pozo 1 en situación normal y en un episodio de inestabilidad de flujo.

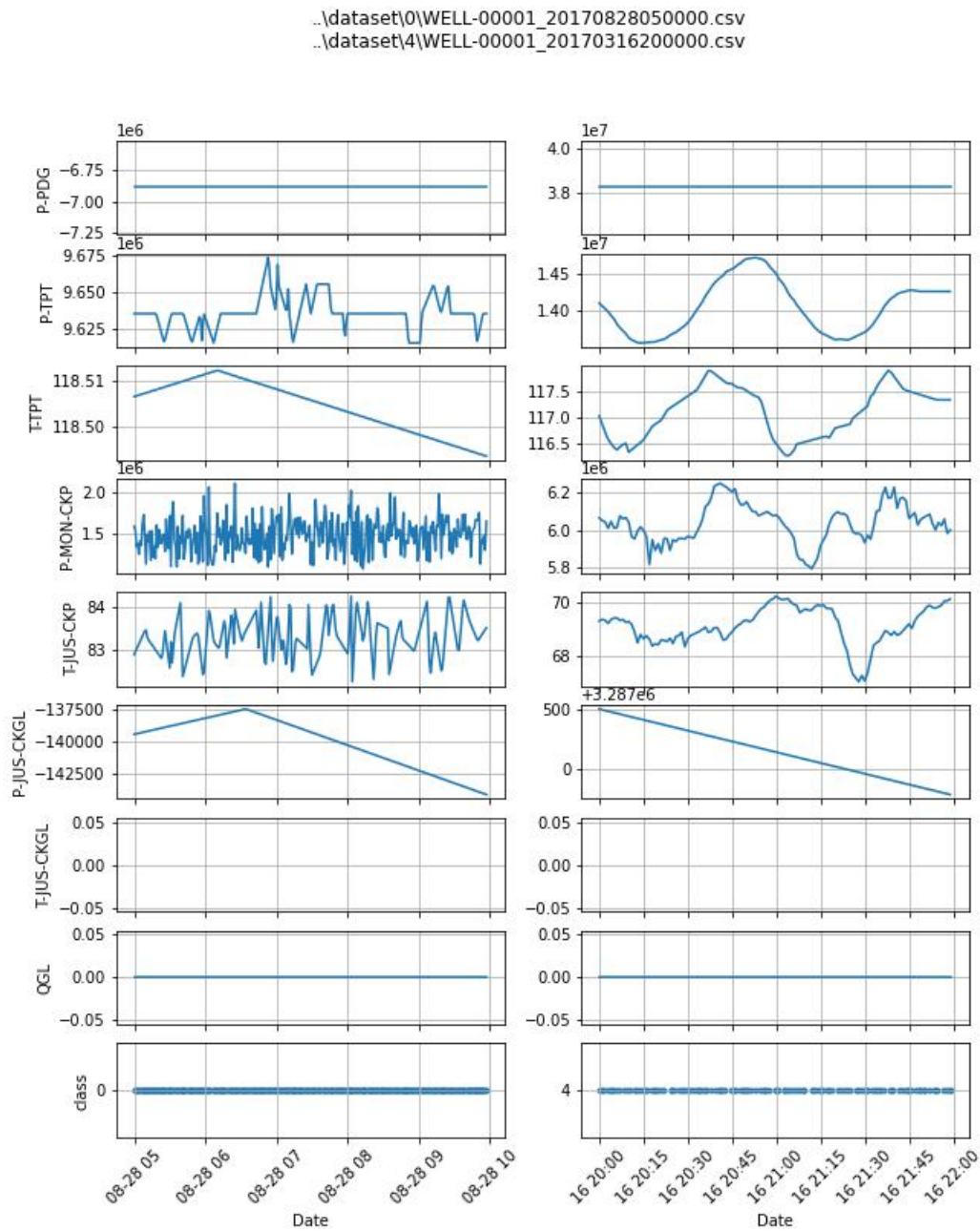


Figura 3.2: Ejemplo de series temporales del dataset 3W

TYPE OF UNDESIRABLE EVENT	WINDOW SIZE
1 - ABRUPT INCREASE OF BSW	12h
2 – SPURIOUS CLOSURE OF DHSV	5min–20min
3 – SEVERE SLUGGING	5h
4 – FLOW INSTABILITY	15min
5 – RAPID PRODUCTIVITY LOSS	12h
6 – QUICK RESTRICTION IN PCK	15min
7 – SCALING IN PCK	72h
8 – HYDRATE IN PRODUCTION LINE	30min–5h

Cuadro 3.1: *Tamaño estimado de las ventanas de tiempo utilizadas para confirmar la ocurrencia de eventos no deseados*[34].

SOURCE	REAL	SIMULATED	HAND-DRAWN	TOTAL
INSTANCE LABEL				
0 - Normal Operation	597	0	0	597
1 - Abrupt Increase of BSW	5	114	10	129
2 - Spurious Closure of DHSV	22	16	0	38
3 - Severe Slugging	32	74	0	106
4 - Flow Instability	344	0	0	344
5 - Rapid Productivity Loss	12	439	0	451
6 - Quick Restriction in PCK	6	215	0	221
7 - Scaling in PCK	4	0	10	14
8 - Hydrate in Production Line	0	81	0	81
TOTAL	1022	939	20	1981

Figura 3.3: Distribución de las series temporales del dataset 3W

En la figura 3.3 se muestra la distribución de series reales(R), simuladas(S) y manuales(D) por tipo de evento.

En la figura 3.4 se muestra el cronograma de adquisición de los datos reales por tipo de evento y pozo.

En el cuadro 3.2 se muestra el porcentaje de valores no disponibles en las series temporales reales del dataset.

label	rcount	P-PDG	P-TPT	T-TPT	P-MON-CKP	T-JUS-CKP	P-JUS-CKGL	T-JUS-CKGL	QGL	class
0	9956791	0.04	0.04	0.04	10.13	14.79	25.63	100.00	25.48	0.00
1	118294	0.00	0.00	0.00	0.00	0.00	0.00	100.00	0.00	0.86
2	158680	0.39	0.39	0.39	69.85	86.66	90.01	100.00	86.66	0.65
3	569152	0.10	0.10	0.10	0.12	0.10	0.12	100.00	0.10	0.00
4	2462076	0.03	0.04	0.04	0.03	0.05	50.51	100.00	23.54	0.00
5	361998	0.00	0.00	0.00	0.00	0.00	0.00	100.00	0.00	0.40
6	54212	0.00	0.00	0.00	0.00	0.00	11.70	100.00	11.70	1.15
7	271708	0.13	0.13	0.13	0.13	0.13	0.13	100.00	0.13	0.23

Cuadro 3.2: Porcentaje de Nan en el Dataset 3W

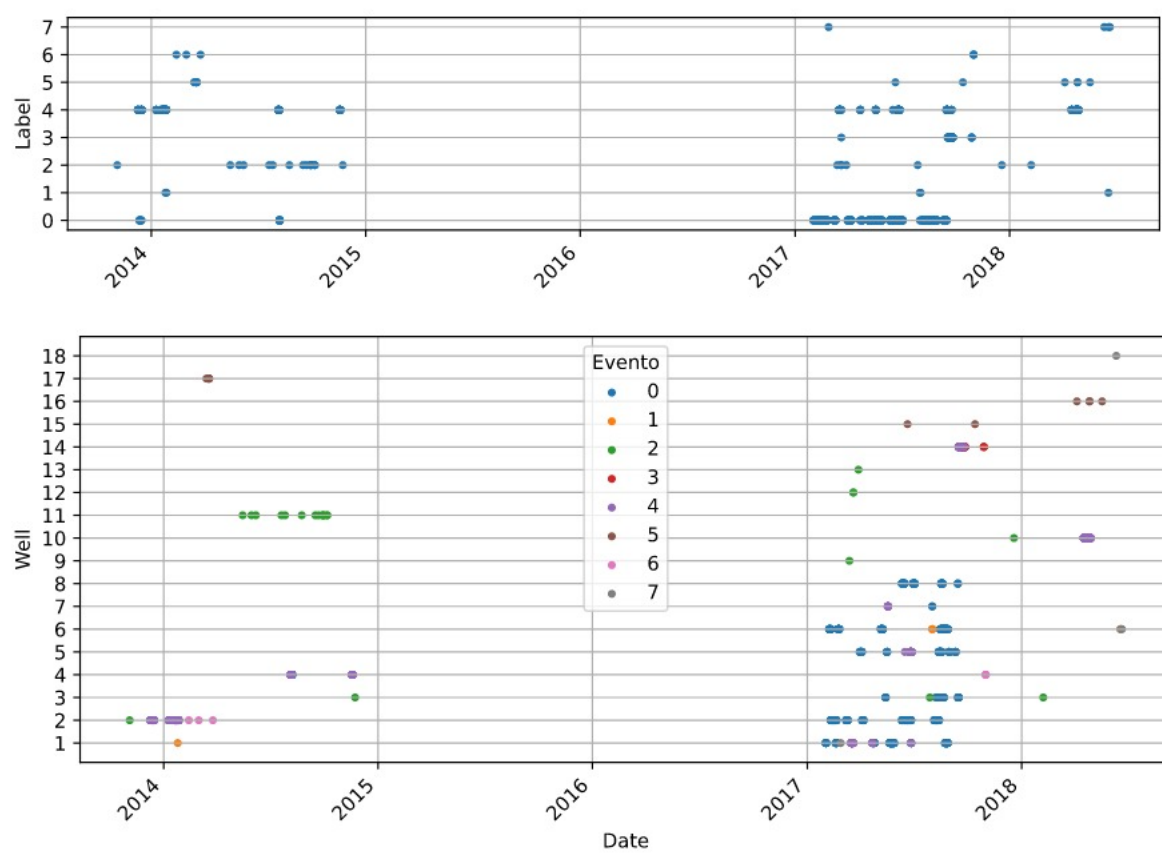


Figura 3.4: Cronograma de la adquisición de datos reales del dataset 3W

3.3. Preprocesamiento de los datos

Las series temporales del Dataset tienen una frecuencia de un registro por segundo. Ya que los eventos que se pretende detectar requieren de unas ventanas de tiempo de entre 5min y 72h [34] no parece que la dinámica del proceso requiera de una frecuencia de muestreo tan alta. Se ha decidido reducir el tamaño del dataset utilizando el promedio minotal de las variables y para la etiqueta **class** el valor más frecuente en el minuto.

Las series temporales tienen como origen datos reales de 18 pozos diferentes, series simuladas y otras rellenadas manualmente. Para poder utilizar el dataset en conjunto las variables se han normalizado (*StandardScaler*). Los datos no disponibles (*Nan*) se han estimado con el método *forward fill* cuando en una serie hay menos del 20 %; si hay más, toda la serie se sustituye por cero. No se ha aplicado ninguna técnica de eliminación de *outliers* al considerar que interferiría con el objetivo de detectar anomalías y drift. En el caso de datos puntuales extraños el promedio de los datos por segundo a minuto los absorberá en parte. En el análisis visual, aunque no ha sido completo, no se han detectado.

Para gestionar los ficheros de las series temporales se ha programado la *python class* **d3w** que permite seleccionar las series según su origen, evento o pozo y generar conjuntos para entrenamiento, test y validación.

Las series seleccionadas con **d3w** se alimenta a una *python class* que genera los datos de forma continua concatenando los registros de cada serie temporal para cada tipo de modelo.

Todo el proyecto se ha desarrollado en *Jupyter Notebooks* y se puede encontrar en https://github.com/BaltiBoix/3W_TFM. Una *fork* del *github* original del Dataset 3W.

3.4. Métodos Offline

Aunque no es el objetivo principal del trabajo se han ajustado diferentes modelos offline al dataset. A continuación se presenta un breve resumen.

3.4.1. Redes Neuronales

Se han ensayado tres modelos de redes neuronales en el entorno Keras relativamente sencillos que se muestran en la figura 3.5. Se ha optado por modelos multiclase utilizando las series temporales de todos los orígenes (reales, simuladas y manuales).

Se ha programado la *python class* **kmodel** que permite definir externamente un modelo

Model: "dnn_model"			Model: "cnn_model"					
Layer (type)	Output Shape	Param #	Layer (type)	Output Shape	Param #			
flatten (Flatten)	(64, 480)	0	layer_normalization (LayerN ormalization)	(64, 30, 16)	32			
layer_normalization (LayerN ormalization)	(64, 480)	960	conv1 (Conv1D)	(64, 30, 128)	4224			
dense (Dense)	(64, 256)	123136	MP1 (MaxPooling1D)	(64, 15, 128)	0			
dropout (Dropout)	(64, 256)	0	drop1 (Dropout)	(64, 15, 128)	0			
dense_1 (Dense)	(64, 128)	32896	conv2 (Conv1D)	(64, 15, 128)	32896			
dropout_1 (Dropout)	(64, 128)	0	MP2 (MaxPooling1D)	(64, 7, 128)	0			
dense_2 (Dense)	(64, 64)	8256	drop2 (Dropout)	(64, 7, 128)	0			
dropout_2 (Dropout)	(64, 64)	0	conv3 (Conv1D)	(64, 7, 128)	32896			
dense_3 (Dense)	(64, 17)	1105	MP3 (MaxPooling1D)	(64, 3, 128)	0			
softmax (Softmax)	(64, 17)	0	drop3 (Dropout)	(64, 3, 128)	0			
			flatten (Flatten)	(64, 384)	0			
			dense (Dense)	(64, 17)	6545			
			softmax (Softmax)	(64, 17)	0			
Total params: 166,353								
Trainable params: 166,353								
Non-trainable params: 0								

Model: "rnn_model"		
Layer (type)	Output Shape	Param #
lstm (LSTM)	(64, 128)	74240
dense1 (Dense)	(64, 17)	2193
softmax (Softmax)	(64, 17)	0
Total params: 76,433		
Trainable params: 76,433		
Non-trainable params: 0		

Figura 3.5: Redes neuronales ensayadas

keras y ensayarlo con el dataset. Se ha utilizado un tamaño de batch de 64 y una ventana de tiempo de 30 minutos. El *learning rate* se ha adaptado a cada tipo de modelo. No se ha hecho una optimización formal de los hiperparámetros por los largos tiempos de ejecución.

Los mejores resultados se obtienen con la RNN y la CNN.

En la figura 3.6 podemos ver la evolución del aprendizaje para la RNN y en la figura 3.7 los resultados para el conjunto de test.

3.4.2. Modelos de Sklearn de aprendizaje incremental

Se han ensayado modelos de Sklearn que permiten el aprendizaje incremental. Eso es, que no requieren cargar todos los datos de una vez y evitan saturar la memoria. Se han probado los modelos *SGDClassifier* y *PassiveAggressiveClassifier* con resultados decepcionantes.

Una mejor solución ha sido primero preprocesar el conjunto de datos a través de un modelo ipca (PCA incremental) y comprimir los datos en conjuntos de datos mucho más pequeños que caben en la memoria y permiten el uso de una mayor variedad de modelos. Utilizando los primeros 24 componentes se captura el 75 % de la varianza. Ajustando un Random Forest se obtienen buenos resultados (Figura 3.8).

La mayoría de las publicaciones sobre el dataset 3W de la bibliografía utilizan modelos offline y obtienen buenos resultados. Por tanto los resultados en este sección del trabajo únicamente confirman que el uso de redes neuronales y Random Forest permite un buen ajuste.



Figura 3.6: Evolución del aprendizaje de la RNN

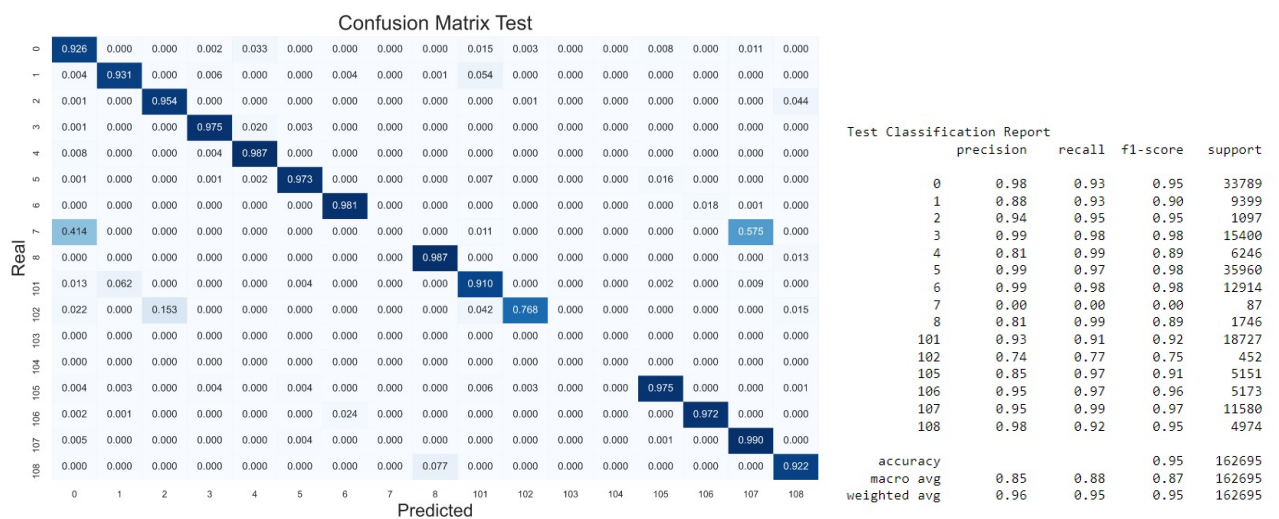


Figura 3.7: Resultados de la RNN para el conjunto de test

	precision	recall	f1-score	support
0	0.81	0.92	0.87	33789
1	0.87	0.84	0.86	9399
2	0.99	0.95	0.97	1097
3	0.97	0.91	0.94	15400
4	0.67	0.77	0.72	6246
5	1.00	0.97	0.98	35960
6	0.99	0.99	0.99	12914
7	1.00	0.83	0.91	87
8	0.84	0.96	0.90	1746
101	0.88	0.86	0.87	18727
102	0.99	0.56	0.71	452
105	0.91	0.97	0.94	5151
106	0.97	0.96	0.97	5173
107	0.95	0.76	0.84	11580
108	0.98	0.87	0.92	4974
accuracy			0.91	162695
macro avg	0.92	0.87	0.89	162695
weighted avg	0.91	0.91	0.91	162695

Figura 3.8: Resultado de aplicar ipca+RF al conjunto de test

3.5. Métodos Online

En esta sección se utiliza el dataset 3W para simular un entorno en tiempo real en el que los datos se procesan registro a registro. Es un ejercicio teórico en el que se hacen algunas aproximaciones que nos alejan del problema real:

- Se considera que se dispone de la etiqueta del registro con solo un minuto de retraso mientras que en el caso real el etiquetado lo han realizado expertos a posteriori.
- Se concatenan las series temporales como si fuera una sola cuando pueden distar entre ellas años (ver el cronograma en la figura 3.4).
- Las series que tienen como origen una simulación o están creadas manualmente no las podemos situar en el tiempo.

3.5.1. Modelo base

Se ha hecho un uso extensivo del módulo River [27]. Para el preprocesamiento de los registros se han utilizado las rutinas que normalizan y calculan variables estadísticas incrementalmente y que, por tanto, no usan información del futuro para el calculo actual. Como en el caso de los métodos offline utilizamos la información de una ventana de tiempo deslizante (*rolling time window*) pero en vez de usar los datos directamente se simplifican en la media, desviación estándar, máximo y mínimo. El uso de datos instantáneos (cada segundo) sería lo que se asemejaría más a recibir los datos directamente de un SCADA pero resulta extremadamente

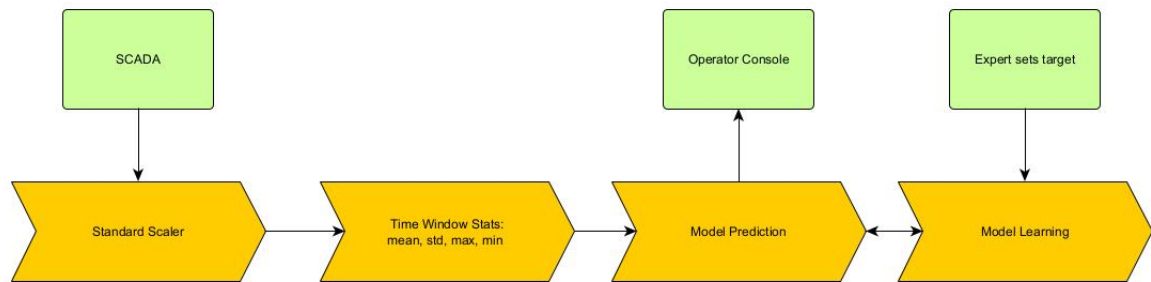
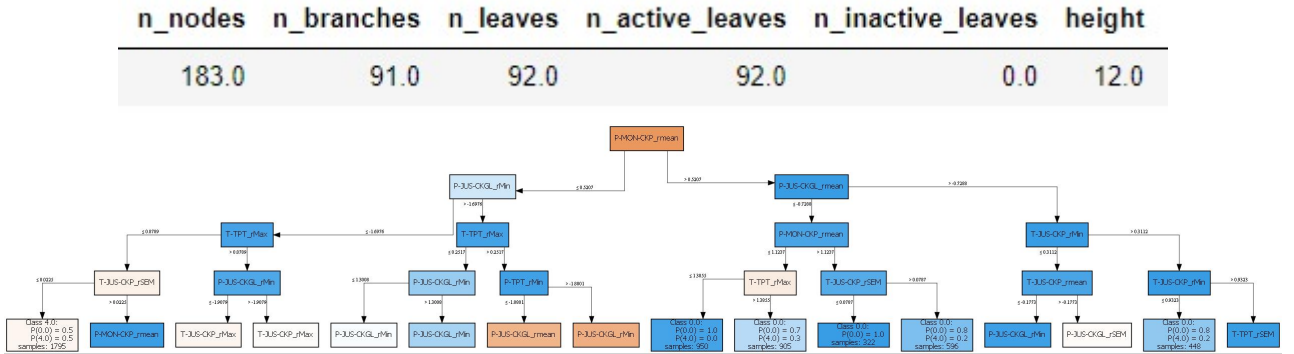


Figura 3.9: Diagrama del modelo online utilizado

lento y se ha optado por usar medias por minuto calculadas con el módulo Pandas. Cada vez que cambia la serie temporal se reinician las variables y no se utilizan en el aprendizaje durante un periodo equivalente a una ventana de tiempo. Las variables preprocesadas se alimentan a un modelo que predice una etiqueta que podemos comparar con la etiqueta que determina el experto y, paso seguido, utilizar para actualizar el modelo según se muestra en el esquema 3.9.

Al evaluar las situaciones anómalas una a una se utilizan las series etiquetadas con esa anomalía junto con las correspondientes a una operación normal. El dataset resultante está desbalanceado a favor de esta última. Para corregirlo se utiliza la rutina de River imblearn RandomOverSampler como *wrapper* para clasificadores. Entrenará al clasificador sobremuestreando el flujo de observaciones para que la distribución de clases vista siga una distribución deseada dada. La implementación es una versión discreta del muestreo de rechazo inverso. El modelo escogido ha sido el clasificador multiclase `tree.HoeffdingTreeClassifier`:

- Un *Hoeffding Tree* es un algoritmo de inducción de árbol de decisión incremental que es capaz de aprender de flujos de datos masivos, suponiendo que los ejemplos de generación de distribución no cambien con el tiempo. Los *Hoeffding Tree* aprovechan el hecho de que una pequeña muestra a menudo puede ser suficiente para elegir un atributo de división óptimo. Esta idea se apoya matemáticamente en el límite de Hoeffding, que cuantifica el número de observaciones (en nuestro caso, ejemplos) necesarias para estimar algunas variables estadísticas dentro de una precisión dada (en nuestro caso, la bondad de un atributo). Una característica teóricamente atractiva de *Hoeffding Tree* que no comparten otros árboles de decisiones incrementales es que tiene un rendimiento garantizado. Usando el límite de Hoeffding, se puede mostrar que su salida es asintóticamente casi idéntica a la de un clasificador no incremental usando infinitos ejemplos. La implementación en River[27] está basada en MOA[5][20].
- Aunque la ejecución del modelo se ha realizado para cada anomalía independientemente se requiere un clasificador multiclase ya que algunas anomalías distinguen tres situaciones (0: Operación Normal, 101: transición, 1: Anomalía).

Figura 3.10: Modelo Hoeffding Tree para la detección de *Flow Instability* con datos reales

BalancedAccuracy: 93.75%					BalancedAccuracy: 88.17%				
	Precision	Recall	F1	Support		Precision	Recall	F1	Support
0.0	98.35%	95.07%	96.68%	149201	0.0	96.45%	92.44%	94.40%	149201
4.0	79.82%	92.43%	85.66%	31501	4.0	70.07%	83.90%	76.37%	31501
Macro	89.08%	93.75%	91.17%		Macro	83.26%	88.17%	85.38%	
Micro	94.61%	94.61%	94.61%		Micro	90.95%	90.95%	90.95%	
Weighted	95.12%	94.61%	94.76%		Weighted	91.85%	90.95%	91.26%	
94.61% accuracy									

Figura 3.11: HT Train Accuracy

Figura 3.12: HT Evaluation Accuracy

Figura 3.13: Precisión del Hoeffding Tree

En la figura 3.10 se muestra una representación gráfica de los primeros 4 de 12 niveles del *Hoeffding Tree* resultante del procesamiento de las series de datos reales con etiquetas 0 (Operación Normal) y 4 (Flujo inestable).

River permite usar modelos de clasificación binarios en entornos multiclase con el uso de los *wrappers* `multiclass.OneVsOneClassifier` y `multiclass.OneVsRestClassifier`. Se ha comparado la precisión obtenida para cada modelo ajustado en modo sólo predicción. Se comportan peor que el Hoeffding Tree. En las figuras 3.13 y 3.16 se puede observar la precisión obtenida en el proceso de aprendizaje incremental y en la evaluación del mismo dataset usando el modelo ya entrenado para el Hoeffding Tree y el `linear_model.ALMAClassifier`.

Los episodios de funcionamiento anormal del proceso de extracción de crudo, cuando se producen, no son puntuales si no que tienen una evolución que se alarga en el tiempo. Este comportamiento tiene como consecuencia que las etiquetas muestren una alta autocorrelación que, por otro lado, es muy frecuente en las series temporales. En la figura 3.17 se compara la precisión del modelo con respecto al modelo de River *dummy.NoChangeClassifier* que simplemente asigna la etiqueta anterior al registro actual. En la parte superior se muestra el cronograma de las series temporales reales utilizadas para la anomalía 4 (*Flow Instability*. Re-

BalancedAccuracy: 99.36%

	Precision	Recall	F1	Support
0.0	99.79%	99.68%	99.74%	149201
4.0	98.50%	99.03%	98.76%	31501
Macro	99.15%	99.36%	99.25%	
Micro	99.57%	99.57%	99.57%	
Weighted	99.57%	99.57%	99.57%	

Figura 3.14: ALMA Train Accuracy

BalancedAccuracy: 59.34%

	Precision	Recall	F1	Support
0.0	90.61%	36.70%	52.24%	149201
4.0	21.47%	81.98%	34.03%	31501
Macro	56.04%	59.34%	43.14%	
Micro	44.60%	44.60%	44.60%	
Weighted	78.56%	44.60%	49.07%	

Figura 3.15: ALMA Evaluation Accuracy

Figura 3.16: Precisión del `linear_model.ALMAClassifier`

pitando el ejercicio para la anomalía 5 (*Rapid Productivity Loss* en el que la mayor parte de las series son simuladas y se han ordenado aleatoriamente se observa el mismo fenómeno ya que la autocorrelación se da dentro de cada serie (Figura 3.18).

3.5.2. Detección de *drift*

El siguiente paso en este trabajo ha sido introducir un detector de *drift* en el esquema 3.9 de procesamiento del dataset visto anteriormente. Cuando se detecta *drift* se reinicia el aprendizaje del modelo siempre que el modelo haya visto un mínimo de registros. Se busca observar si se obtiene una mejora de precisión global. El esquema modificado se muestra en la figura 3.19.

El detector de *drift* recibe la etiqueta predicha y la real del experto y estima si se ha producido un cambio en la distribución del error. El módulo River cuenta con varios métodos de detección que se han evaluado con el dataset 3W:

- **ADWIN** (ADaptive WINdowing) es un método popular de detección de *drift*. ADWIN mantiene una ventana de longitud variable de elementos recientes. Esta ventana se divide a su vez en dos subventanas (W_0 , W_1) que se utilizan para determinar si se ha producido un cambio en la distribución de los datos. ADWIN compara el promedio de W_0 y W_1 para confirmar que corresponden a la misma distribución. El *concept drift* se detecta si la igualdad de distribución ya no se mantiene. Al detectar un desvío, W_0 se reemplaza por W_1 y se inicializa uno nuevo W_1 . ADWIN utiliza un valor significativo $\delta \in (0, 1)$ para determinar si las dos subventanas corresponden a la misma distribución [4].
- **DDM** (Método de detección de deriva) es un método de detección de *concept drift* basado en la premisa del modelo de aprendizaje PAC, que la tasa de error del clasificador disminuirá a medida que aumente la cantidad de muestras analizadas, siempre que la distribución de datos sea estacionaria.

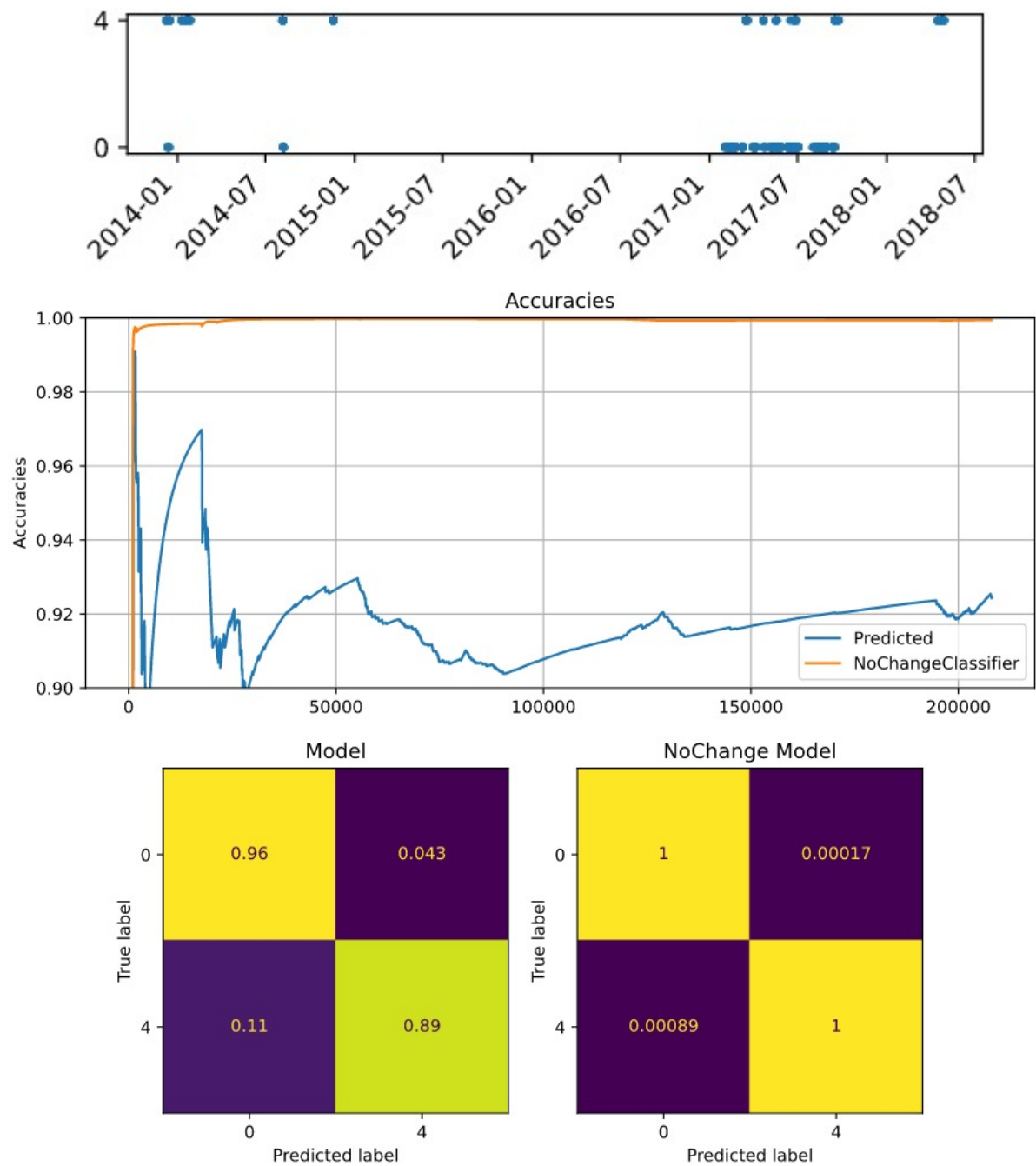


Figura 3.17: Comparación del modelo Hoeffding Tree con el NoChange para *Flow Instability* con datos reales

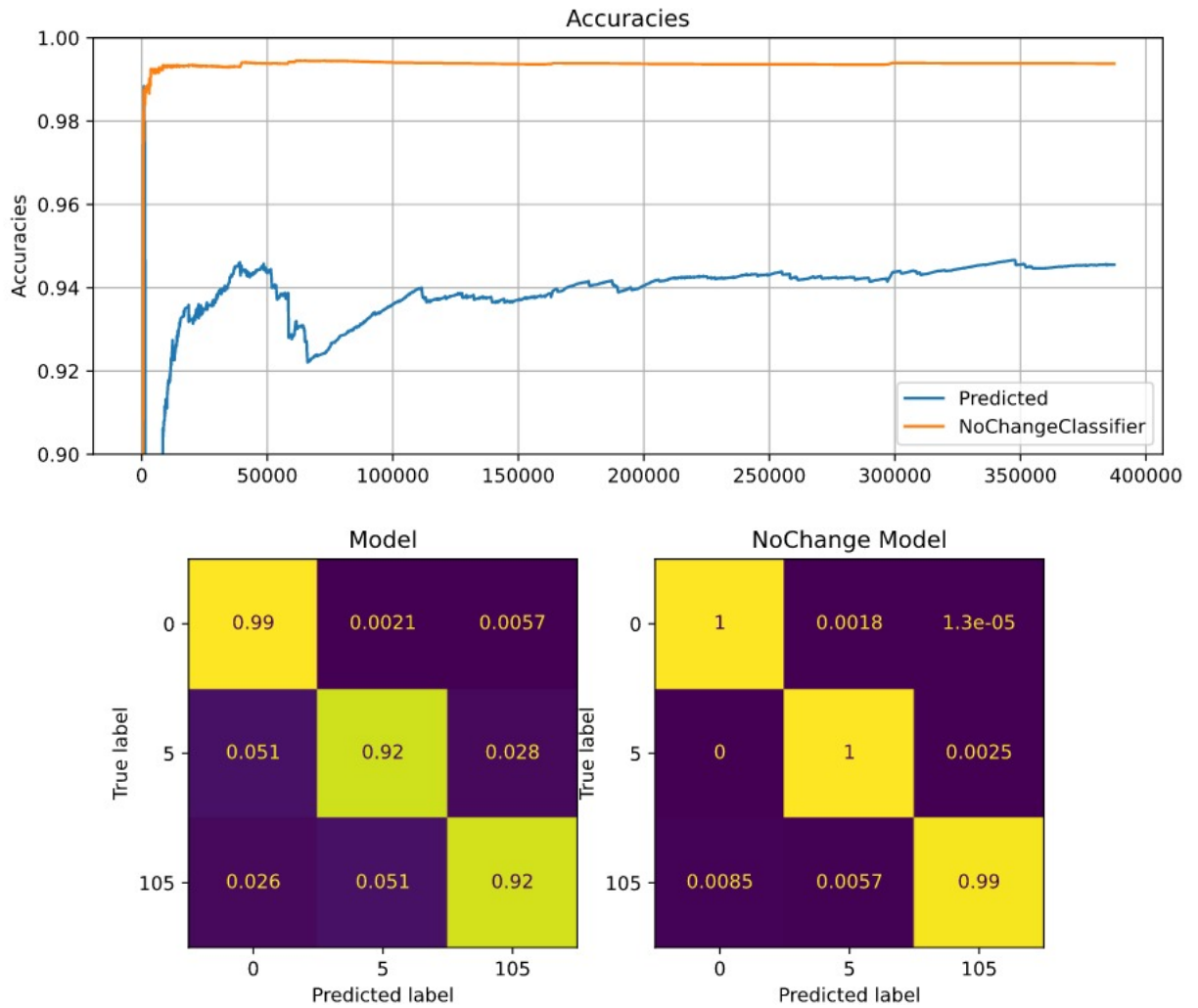


Figura 3.18: Comparación del modelo Hoeffding Tree con el NoChange para *Rapid Productivity Loss* con datos reales y simulados ordenados aleatoriamente

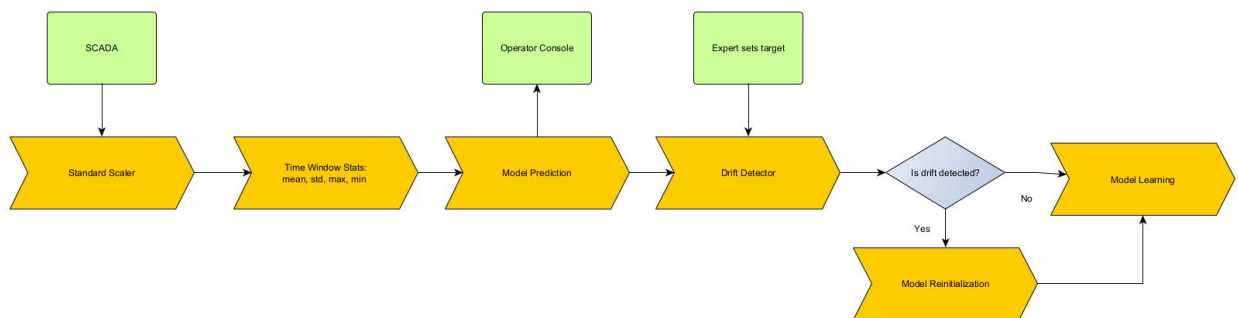


Figura 3.19: Diagrama del modelo online utilizado con detección de *drift*

Si el algoritmo detecta un aumento en la tasa de error, que supera un umbral calculado, se detecta un cambio o el algoritmo advertirá al usuario que puede ocurrir un cambio en un futuro próximo, lo que se denomina zona de advertencia (*Warning Zone*).

El umbral de detección se calcula en función de dos estadísticas, obtenidas cuando es $(p_i + s_i)$ mínimo:

- p_{min} : La tasa de error mínima registrada.
- s_{min} : La desviación estándar mínima registrada.

En el instante i , el algoritmo de detección utiliza:

- p_i : La tasa de error en el instante i .
- s_i : La desviación estándar en el instante i .

Las condiciones para ingresar a la zona de advertencia y detectar cambios son las siguientes:

- si $p_i + s_i \geq p_{min} + w_l * s_{min} \Rightarrow$ Zona de advertencia
- si $p_i + s_i \geq p_{min} + d_l * s_{min} \Rightarrow$ Cambio detectado

En las expresiones anteriores, w_l y d_l representan, respectivamente, los umbrales de advertencia y deriva (*drift*) [14].

- **EDDM** (*Early Drift Detection Method*) tiene como objetivo mejorar la tasa de detección del *concept drift* gradual en DDM, manteniendo un buen rendimiento en la abrupta.

Este método funciona al realizar un seguimiento de la distancia promedio entre dos errores en lugar de solo la tasa de error. Para esto, es necesario realizar un seguimiento de la distancia promedio y la desviación estándar, así como la distancia máxima y la desviación estándar máxima.

El algoritmo funciona de manera similar al algoritmo DDM, al realizar un seguimiento de las estadísticas únicamente. Funciona con la distancia media (p'_i) y la desviación estándar (s'_i), así como p'_{max} y s'_{max} , que son los valores de p'_i y s'_i cuando $(p'_i + 2 * s'_i)$ alcanza su máximo.

Al igual que DDM, hay dos valores de umbral que definen el límite entre no cambio, zona de advertencia y desviación detectada:

- si $(p'_i + 2 * s'_i) / (p'_{max} + 2 * s'_{max}) < \alpha \Rightarrow$ Zona de advertencia
- si $(p'_i + 2 * s'_i) / (p'_{max} + 2 * s'_{max}) < \beta \Rightarrow$ Cambio detectado

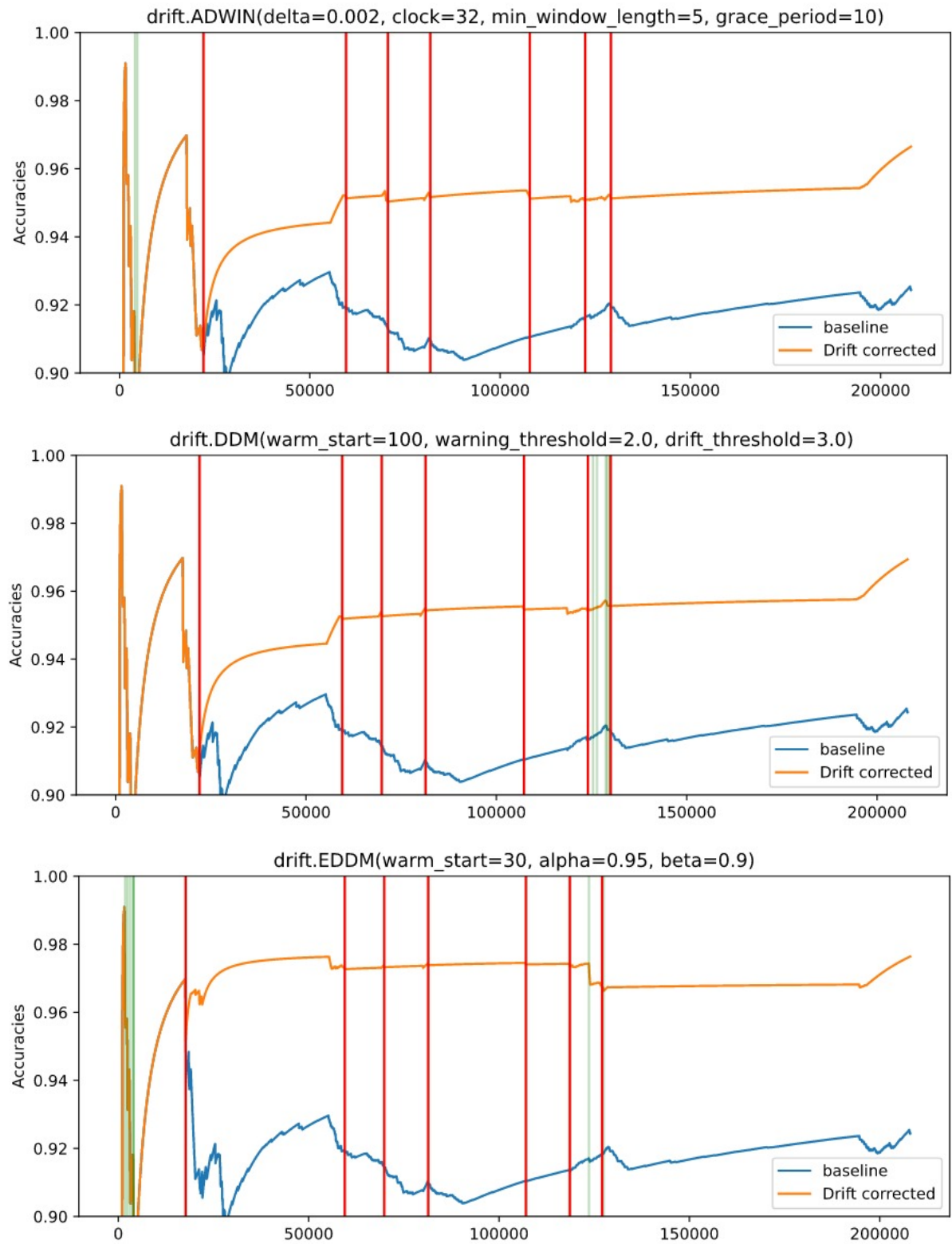
α y β por defecto valen 0.95 y 0.9, respectivamente [3].

- **HDDM_A** es un método de detección de *drift* basado en la desigualdad de Hoeffding que utiliza el promedio de entrada como estimador[5][13].
- **HDDM_W** es un método de detección de *drift* en línea basado en los límites de McDiarmid. Utiliza la media móvil ponderada exponencialmente (EWMA) como estimador [5][13].

En los gráficos de la figura 3.20 se puede ver la aplicación de los métodos de detección de *drift* descritos anteriormente a las series de datos reales con etiquetas 0 (Operación Normal) y 4 (Flujo inestable) ordenados por fecha. Se reinicializa el aprendizaje del modelo cuando se detecta *drift* y el modelo al menos ha visto 5000 registros. Se observa que la mayor parte de las detecciones ocurren en aproximadamente los mismos puntos y en todos los casos mejora la precisión.

El comportamiento del resto de anomalías no es significativo ya que como se muestra en la figura 3.3 se dispone de pocas series reales. se muestran los resultados para las anomalías 1 a 3 (ver cuadro 3.1) con el método de detección ADWIN en la figura 3.21.

Repitiendo el ejercicio para la anomalía 5 (*Rapid Productivity Loss* con series simuladas y ordenadas las series aleatoriamente (no los registros) con los mismos ajustes de ADWIN que se han utilizado anteriormente se obtiene un resultado absurdo (Figura 3.22).



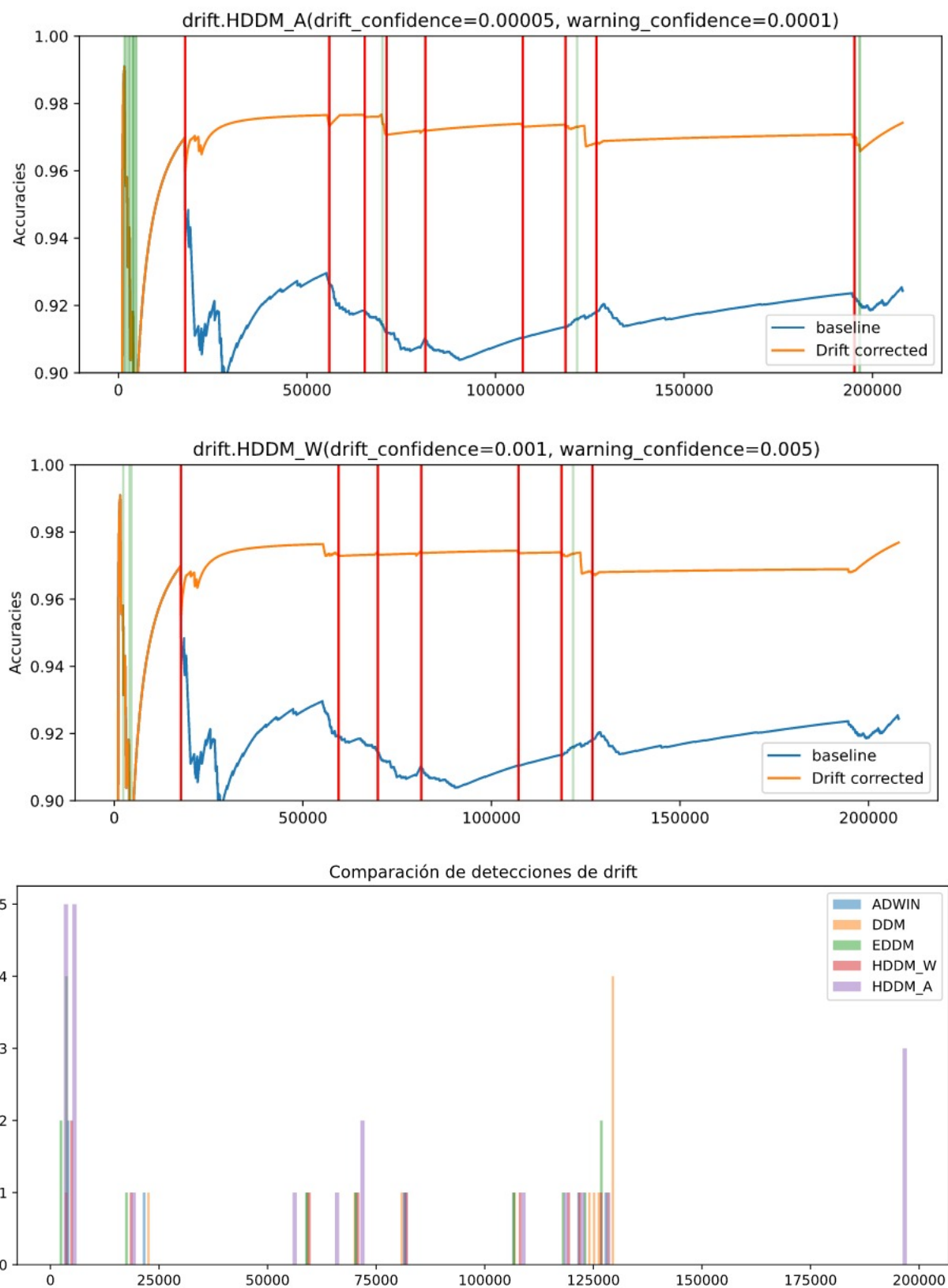


Figura 3.20: Detección de *drift* con diferentes métodos para *flow instability*

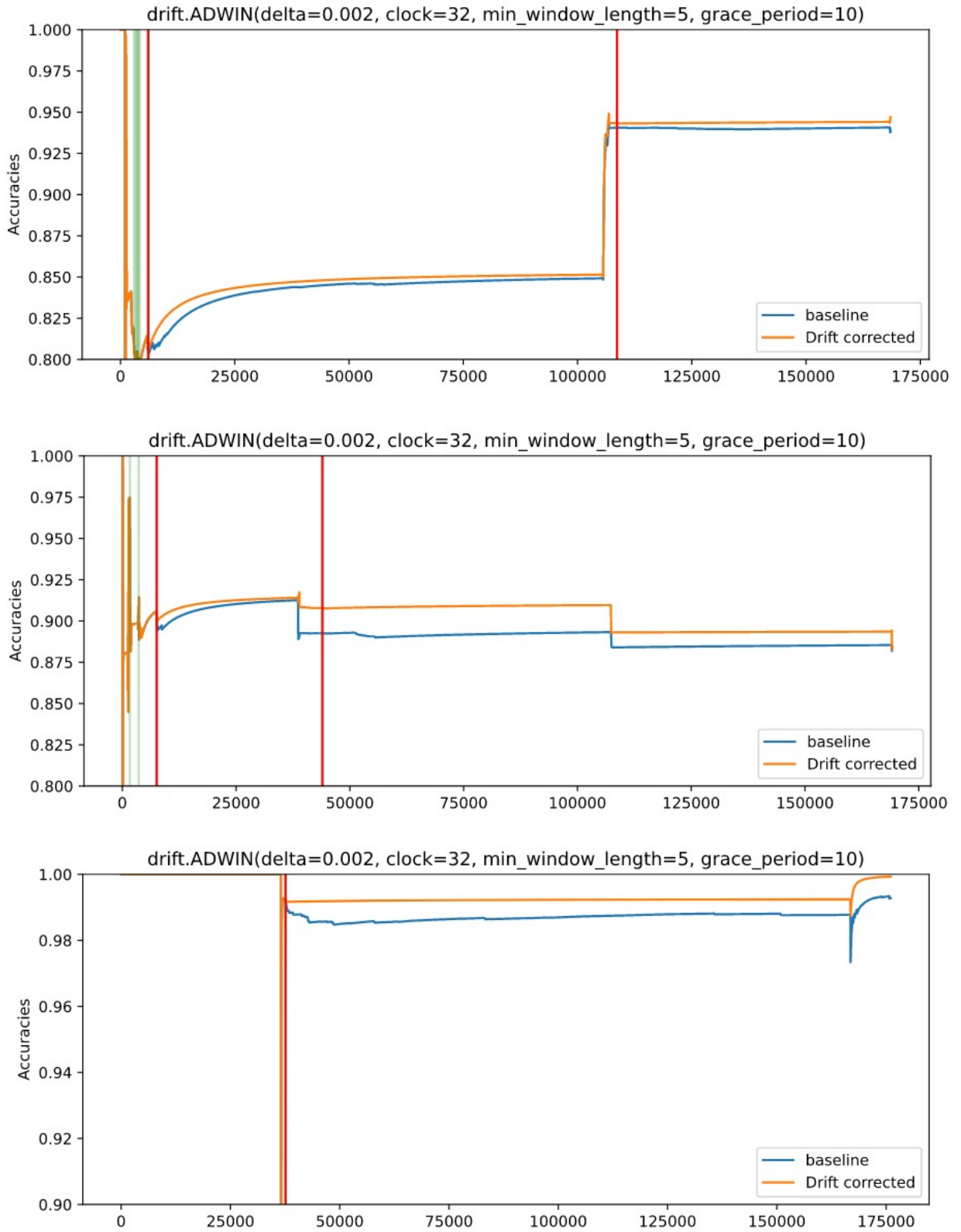


Figura 3.21: Detección de *drift* con ADWIN para anomalías 1 (ABRUPT INCREASE OF BSW), 2 (SPURIOUS CLOSURE OF DHSV) y 3 (SEVERE SLUGGING)

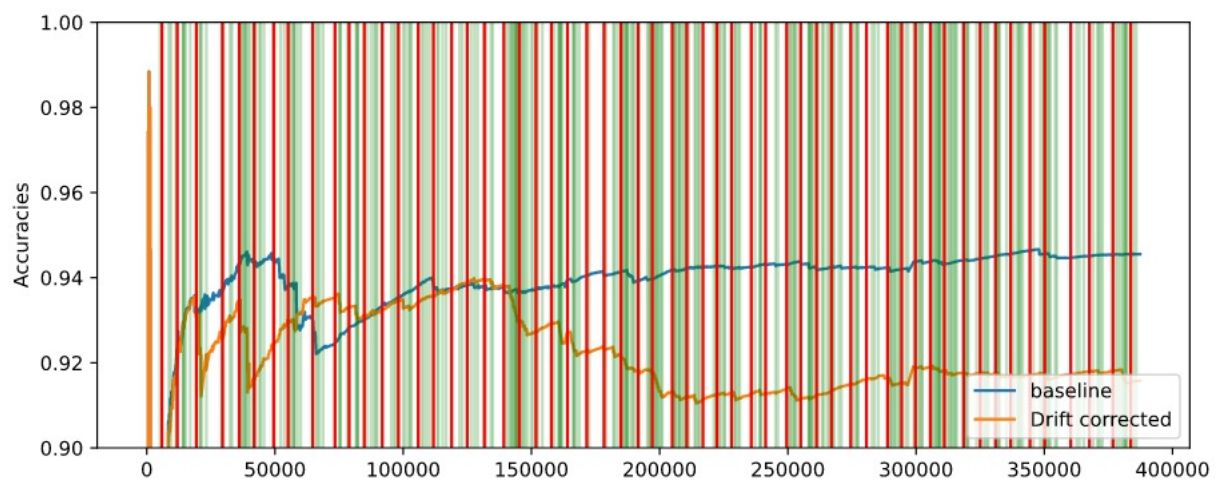


Figura 3.22: Detección de *drift* con ADWIN para la anomalía 5 (*Rapid Productivity Loss*) con series simuladas

Capítulo 4

Conclusiones del TFM

De los estudios publicados hasta la fecha y de este trabajo podemos concluir que es posible desarrollar una aplicación de soporte a la operación de los pozos de petróleo *offshore* que explote los modelos offline ajustados con el dataset 3W. En este trabajo se han obtenido buenos resultados utilizando *Random Forest* y redes neuronal RNN y CNN relativamente sencillas.

A partir de las series temporales del dataset 3W se ha emulado la recepción de los datos desde un SCADA (Supervisory Control and Data Acquisition) y su procesamiento utilizando las rutinas de *stream learning* del módulo River[27] de preprocesamiento y aprendizaje incremental. A este proceso se le ha implementado un detector de *concept drift* que reinicializa el aprendizaje. Se observa una clara mejora cuando se utilizan datos reales de dataset ordenados por fecha. Lamentablemente únicamente en el caso de la anomalía *Flow Instability* se dispone de suficientes datos reales (no simulados o manuales).

Una posible explicación del *concept drift* detectado podría ser que los expertos que han etiquetado las series hayan cambiado con los años o haya cambiado el criterio aplicado.

Aunque los métodos de detección de *drift* han dado resultados similares, ADWIN es el más configurable y permite un ajuste más fino.

Hay que resaltar que el uso de un proceso de aprendizaje incremental como el que se ha emulado sólo sería posible realmente si dispusiéramos de las etiquetas con un minuto de retraso. En realidad el etiquetado lo realiza un experto offline. En cualquier caso sería válido para modelos cuyo objetivo sea hacer predicciones (*forecasting*) ya que se evita el etiquetado offline.

Un posterior trabajo podría ser el uso de un modelo ajustado offline con datos de los años 2014 y 2015 con un detector de *drift* aplicado sobre los datos de entrada al modelo (en vez de sobre el error en la predicción)[31] que reportara si la predicción deja de ser fiable con datos de los años posteriores.

Bibliografía

- [1] Charu C. Aggarwal. *Outlier Analysis*. Springer, 2013.
- [2] Nida Aslam, Irfan Ullah Khan, Aisha Alansari, Marah Alrammah, Atheer Alghwairy, Rahaf Alqahtani, Razan Alqahtani, Maryam Almushikes, and Mohammed AL Hashim. Anomaly detection using explainable random forest for the prediction of undesirable events in oil wells. *Applied Computational Intelligence and Soft Computing*, 2022:1558381, Aug 2022.
- [3] Manuel Baena-García, José Campo-Ávila, Raúl Fidalgo-Merino, Albert Bifet, Ricard Gavaldà, and Rafael Morales-Bueno. Early drift detection method. *Fourth International Workshop on Knowledge Discovery from Data Streams*, 01 2006.
- [4] Albert Bifet and Ricard Gavaldà. Learning from time-changing data with adaptive windowing. *Proceedings of the 7th SIAM International Conference on Data Mining*, 7, 04 2007.
- [5] Albert Bifet, Geoffrey Holmes, Richard Kirkby, and Bernhard Pfahringer. Moa: massive online analysis. *Journal of Machine Learning Research*, 11, 05 2010.
- [6] Ekaba Bisong. *Batch vs. Online Learning*, pages 199–201. Apress, Berkeley, CA, 2019.
- [7] Chrisander Brønstad, Sergio L. Netto, and Antonio L.L. Ramos. In *Data-driven Detection and Identification of Undesirable Events in Subsea Oil Wells*, Athene, Greece, 2021. SENSORDEVICES 2021. The Twelfth International Conference on Sensor Device Technologies and Applications.
- [8] Bruno Guilherme Carvalho, Ricardo Emanuel Vaz Vargas, Ricardo Menezes Salgado, Celso José Munaro, and Flávio Miguel Varejão. Hyperparameter tuning and feature selection for improving flow instability detection in offshore oil wells. In *2021 IEEE 19th International Conference on Industrial Informatics (INDIN)*, pages 1–6, 2021.
- [9] Bruno Guilherme Carvalho, Ricardo Emanuel Vaz Vargas, Ricardo Menezes Salgado, Celso José Munaro, and Flávio Miguel Varejão. Flow instability detection in offshore oil wells

- with multivariate time series machine learning classifiers. In *2021 IEEE 30th International Symposium on Industrial Electronics (ISIE)*, pages 1–6, 2021.
- [10] A. De Salvo Castro, M. De Jesus Rocha Santos, F. Leta, C. Lima, and G. Lima. Unsupervised methods to classify real data from offshore wells. *American Journal of Operations Research*, 11:227–241, 2021.
- [11] Agnieszka Duraj and Piotr S. Szczepaniak. Outlier detection in data streams — a comparative study of selected methods. *Procedia Computer Science*, 192:2769–2778, 2021. Knowledge-Based and Intelligent Information & Engineering Systems: Proceedings of the 25th International Conference KES2021.
- [12] Ilan Sousa Figueirêdo, Tássio Farias Carvalho, Wenisten José Dantas Silva, Lílían Lefol Nani Guarieiro, and Erick Giovani Sperandio Nascimento. In *Detecting Interesting and Anomalous Patterns In Multivariate Time-Series Data in an Offshore Platform Using Unsupervised Learning*, volume Day 2 Tue, August 17, 2021 of *OTC Offshore Technology Conference*, 08 2021. D021S025R003.
- [13] Isvani Frías-Blanco, José del Campo-Ávila, Gonzalo Ramos-Jiménez, Rafael Morales-Bueno, Agustín Ortiz-Díaz, and Yailé Caballero-Mota. Online and non-parametric drift detection methods based on hoeffding’s bounds. *IEEE Transactions on Knowledge and Data Engineering*, 27(3):810–823, 2015.
- [14] João Gama, Pedro Medas, Gladys Castillo, and Pedro Pereira Rodrigues. Learning with drift detection. In Ana L. C. Bazzan and Sofiane Labidi, editors, *SBIA*, volume 3171 of *Lecture Notes in Computer Science*, pages 286–295. Springer, 2004.
- [15] João Gama, Indrè Žliobaitė, Albert Bifet, Mykola Pechenizkiy, and Abdelhamid Bouchachia. A survey on concept drift adaptation. *ACM Comput. Surv.*, 46(4), mar 2014.
- [16] Rosana Noronha Gemaque, Albert França Josuá Costa, Rafael Giusti, and Eulanda Miranda dos Santos. An overview of unsupervised drift detection methods. *WIREs Data Mining and Knowledge Discovery*, 10(6):e1381, 2020.
- [17] Heitor Murilo Gomes, Jesse Read, Albert Bifet, Jean Paul Barddal, and Joao Gama. Machine learning for streaming data: state of the art, challenges, and opportunities. *ACM SIGKDD Explorations Newsletter*, 21:6–22, 11 2019.
- [18] Manish Gupta, Jing Gao, Charu C. Aggarwal, and Jiawei Han. Outlier detection for temporal data: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 26(9):2250–2267, Sep. 2014.

- [19] Hanqing Hu, Mehmed Kantardzic, and Tegjyot S. Sethi. No free lunch theorem for concept drift detection in streaming data classification: A review. *WIREs Data Mining and Knowledge Discovery*, 10(2):e1327, 2020.
- [20] Geoff Hulten, Laurie Spencer, and Pedro Domingos. Mining time-changing data streams. In *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '01, page 97–106, New York, NY, USA, 2001. Association for Computing Machinery.
- [21] Bartosz Krawczyk, Leandro L. Minku, João Gama, Jerzy Stefanowski, and Michał Woźniak. Ensemble learning for data stream analysis: A survey. *Information Fusion*, 37:132–156, 2017.
- [22] Jesus L. Lobo, Javier Del Ser, Albert Bifet, and Nikola Kasabov. Spiking neural networks and online learning: An overview and perspectives. *Neural Networks*, 121:88–100, 2020.
- [23] Jie Lu, Anjin Liu, Fan Dong, Feng Gu, João Gama, and Guangquan Zhang. Learning under concept drift: A review. *IEEE Transactions on Knowledge and Data Engineering*, PP:1–1, 10 2018.
- [24] Jesús López Lobo. *New perspectives and methods for stream learning in the presence of concept drift*. PhD thesis, UNIVERSITY OF THE BASQUE COUNTRY UPV/EHU, Spain, 2018.
- [25] André Paulo Ferreira Machado, Ricardo Emanuel Vaz Vargas, Patrick Marques Ciarelli, and Celso Jose Munaro. Improving performance of one-class classifiers applied to anomaly detection in oil wells. *Journal of Petroleum Science and Engineering*, 218:110983, 2022.
- [26] Matheus A. Marins, Bettina D. Barros, Ismael H. Santos, Daniel C. Barrionuevo, Ricardo E.V. Vargas, Thiago de M. Prego, Amaro A. de Lima, Marcello L.R. de Campos, Eduardo A.B. da Silva, and Sergio L. Netto. Fault detection and classification in oil wells and production/service lines using random forest. *Journal of Petroleum Science and Engineering*, 197:107879, 2021.
- [27] Jacob Montiel, Max Halford, Saulo Martiello Mastelini, Geoffrey Bolmier, Raphael Sourty, Robin Vaysse, Adil Zouitine, Heitor Murilo Gomes, Jesse Read, Talel Abdessalem, and Albert Bifet. *River: machine learning for streaming data in python*, 2020.
- [28] National Commission on the BP Deepwater Horizon Oil Spill and Offshore Drilling. *Deep water : the Gulf oil disaster and the future of offshore drilling : report to the President*. Washington, D.C, 2011.

-
- [29] Dragoljub Pokrajac, Aleksandar Lazarevic, and Longin Jan Latecki. Incremental local outlier detection for data streams. *2007 IEEE Symposium on Computational Intelligence and Data Mining*, pages 504–515, 2007.
- [30] James Reason. Human error: models and management. *BMJ*, 320(7237):768–770, 2000.
- [31] Tegjyot Singh Sethi and Mehmed Kantardzic. On the reliable detection of concept drift from streaming unlabeled data. *Expert Systems with Applications*, 82:77–99, 2017.
- [32] Alexey Tsymbal. The problem of concept drift: Definitions and related work. 05 2004.
- [33] Evren Mert Turan and Johannes Jäschke. Classification of undesirable events in oil well operation. *2021 23rd International Conference on Process Control (PC)*, pages 157–162, 2021.
- [34] Ricardo Emanuel Vaz Vargas, Celso José Munaro, Patrick Marques Ciarelli, André Gonçalves Medeiros, Bruno Guberfain do Amaral, Daniel Centurion Barrionuevo, Jean Carlos Dias de Araújo, Jorge Lins Ribeiro, and Lucas Pierezan Magalhães. A realistic and public dataset with rare undesirable real events in oil wells. *Journal of Petroleum Science and Engineering*, 181, 7 2019.
- [35] Venkat Venkatasubramanian, Raghunathan Rengaswamy, Kewen Yin, and Surya N. Kavuri. A review of process fault detection and diagnosis: Part i: Quantitative model-based methods. *Computers & Chemical Engineering*, 27(3):293–311, 2003.