

# Presentación

---

## Diapositiva 1 (Introducción)

---

Bienvenidos a la presentación del trabajo fin de master: "Detección de eventos anómalos en un entorno industrial mediante el uso de técnicas de Federated Learning" mi nombre es Darío Martín García Carretero. A lo largo de esta presentación se va a mostrar una de las muchas aplicaciones del uso del aprendizaje federado. En este caso se utilizará como herramienta para la detección de eventos anómalos en un entorno industrial. Sin más preámbulos comencemos con la presentación.

## Diapositiva 2 (Contexto y justificación)

---

### Bloque 1

¿Qué entendemos por un evento anómalo?

Un evento anómalo es aquel que se produce de forma repentina y sin previsión. Estos eventos anómalos pueden ser de muy diversas índoles dependiendo del contexto. Por ejemplo, dentro de la medicina un evento anómalo podría ser una arritmia cardiaca detectada en un paciente o, en el contexto de la seguridad podría ser el sonido del disparo de una pistola que ha sido registrado en las grabaciones de una cámara de seguridad.

### Bloque 2

---

En un entorno industrial, los eventos anómalos también pueden ser de muy diversa índole, en nuestro trabajo consideraremos como evento anómalo los fallos que, irremediablemente se producen en las máquinas debido al desgaste derivado de su uso continuado.

### Bloque 3

---

Hay que tener en cuenta que estos fallos pueden provocar grandes daños económicos y personales. Por lo que su detección es crucial y puede ayudar a prevenir situaciones irreversibles.

## Diapositiva 3 (Escenario I)

---

Pasemos ahora a exponer un caso concreto de posible aplicación del Federated Learning. Este caso de uso será en el que se basará el resto del trabajo.

Supongamos que una empresa desea implantar un sistema que detecte las posibles averías de las máquinas de sus factorías antes de que estas fallen y puedan provocar daños tanto al personal como a las propias instalaciones.

### Bloque 1

Sabemos que empresa dispone de plantas repartidas por todo el mundo. A pesar de pertenecer a la misma compañía cada una de las instalaciones tiene sus propias particularidades en cuanto al tipo de productos que fabrican, la manera de fabricarlos, etc. A estas diferencias se han de añadir también las condiciones ambientales de cada lugar: temperatura, humedad, presión atmosférica, etc. A pesar de estas diferencias, las máquinas utilizadas en todas las factorías son similares.

### Bloque 2

---

Las plantas trabajan con un alto grado de independencia y de hecho, suelen competir entre ellas en cuestiones como: la cantidad

de producción, calidad, etc. Debido a esta competitividad, las factorías son reacias a compartir datos sobre sus técnicas de producción, las configuraciones de sus máquinas, etc. Esto implica que el acceso a sus datos está muy restringido y que estos únicamente pueden ser utilizados a nivel interno.

## Bloque 3

---

La compañía está en constante expansión y es habitual que abra nuevas plantas a lo largo de mundo. Obviamente la empresa desea implantar el sistema de detección de fallos en estas nuevas fábricas lo más rápidamente posible.

## Diapositiva 4 (Escenario II)

---

### Bloque 1

Hoy en día la mayoría de los componentes dentro de un entorno industrial están monitorizados mediante el uso de dispositivos de medición especializados y esta empresa no es una excepción y por lo tanto, es posible disponer de una gran cantidad de datos de los sensores ubicados en las factorías.

### Bloque 2

Como se dispone de esa gran cantidad de datos sobre el funcionamiento de las instalaciones (hay que recordar que los datos solo se pueden usar a nivel interno, es decir, de manera local) es totalmente lógico tratar de resolver el problema mediante el uso de técnicas de Machine Learning.

### Bloque 3

A pesar de poder aplicar técnicas de Machine Learning tradicional para resolver el problema a nivel local, es importante darse cuenta de que es imprescindible tratar el problema a nivel global por el siguiente motivo, la necesidad de una implantación rápida del sistema en las instalaciones de nueva creación. Si abordásemos el problema de forma tradicional tendríamos dos alternativas:

- Repetir el proceso de la creación de modelos que se siguió en el resto de las plantas. Esto podría llevar meses o incluso años, lo que iría en contra del requerimiento de rápida implantación.
- Juntar los datos de todas las plantas en un único lugar y construir un modelo con todos los datos. Esto proporcionaría un modelo más robusto capaz de adaptarse a las condiciones de las nuevas instalaciones. El problema es que esta aproximación violaría el requerimiento de privacidad de los datos.

Por estos motivos se propone el uso del Federated Learning que permite crear modelos de gran calidad y a la vez cumplir con las restricciones de privacidad y rápido despliegue.

Hemos mencionado conceptos como Machine Learning, Federated Learning pero no hemos explicado en qué consisten. En las siguientes diapositivas realizaremos una breve introducción de estos conceptos.

## Diapositiva 5 (Machine Learning I)

---

### Bloque 1

¿Qué es el Machine Learning? El Machine Learning (en español, aprendizaje automático) es un subcampo de la computación y una rama de la inteligencia artificial. Su objetivo es crear programas (llamados comúnmente modelos) capaces de generalizar comportamientos a través de la información suministrada en forma de ejemplos (también llamados instancias).

### Bloque 2

¿Qué puede hacer el Machine Learning por nosotros?, ¿Qué aplicaciones tiene? Hoy en día el Machine Learning tiene una gran

variedad de aplicaciones, entre las que se incluyen: motores de búsqueda, diagnóstico médico, detección de fraude en el uso de tarjetas de crédito, clasificación de secuencias de ADN, videojuegos, etc.

## Diapositiva 6 (Machine Learning II)

---

Bien, ¿Y cómo funciona? Para explicarlo de una manera sencilla estableceremos una analogía entre nuestro cerebro y un modelo de Machine Learning mediante un simple ejemplo:

### Bloque 1

Cuando somos niños y vemos por primera vez un pez no reconoceremos el objeto porque nunca hemos visto algo similar. Si nos explican de que se trata, la próxima vez que veamos un pez, lo reconoceremos inmediatamente. Esto sucede debido a que, de forma inconsciente, nuestro cerebro ha almacenado las características del animal: que tiene aletas, cola, escamas, etc.

### Bloque 2

El aprendizaje automático funciona, en la mayoría de los casos, de forma análoga. Al modelo (nuestro cerebro) se le suministra un conjunto de datos etiquetados (pez o no pez) y el modelo "aprende" a reconocer patrones en los datos. Posteriormente ese modelo, gracias a la generalización, será capaz de reconocer peces cuando los vea.

## Diapositiva 7 (Machine Learning III)

---

¿Cómo se aplica el Machine Learning? La construcción de un modelo de aprendizaje automático puede dividirse, a grandes rasgos, en tres etapas:

### Bloque 1

Etapa 1, adquisición y preparación de los datos. Como ya se mencionó, los modelos "aprenden" mediante el uso de ejemplos etiquetados. Es, en esta etapa, donde se obtienen estos ejemplos. Las fuentes pueden ser variadas y dependen del escenario en el que estemos trabajando. Por ejemplo, en nuestro caso, los datos provendrían de dispositivos de medición ubicados en las factorías y de los registros de mantenimiento de las máquinas. Esta es la parte de la adquisición. Es muy frecuente que los datos capturados no estén en un formato adecuado y deben ser procesados para darles la estructura necesaria para poder entrenar a los modelos esta tarea es vital y se denomina: preparación de los datos.

### Bloque 2

Etapa 2, entrenamiento del modelo. Esta es la etapa donde el modelo "aprende" gracias a los datos adquiridos y preparados en la etapa anterior.

### Bloque 3

Etapa 3, validación del modelo. Esta es la etapa donde se "examina" el modelo construido. Es decir, donde se valora si ha aprendido lo suficiente. Aunque existen multitud de técnicas, la forma más habitual es la de suministrarle al modelo un conjunto de datos que no hayan sido utilizados en el entrenamiento y evaluar sus resultados en función de una o varias métricas.

### Bloque 4

Es importante destacar que, aunque estas fases se han presentado de manera secuencial, no necesariamente se aplican de esa manera y es habitual que se salte de una a otra dependiendo de las características del problema tratado.

## Diapositiva 8 (Federated Learning)

---

Veamos ahora en que consiste el Federated Learning (aprendizaje federado en español). El aprendizaje federado es una técnica de aprendizaje automático que entrena un algoritmo a través de múltiples dispositivos o servidores, denominados nodos sin intercambiar datos entre ellos. El proceso de entrenamiento se divide en 4 fases:

## Bloque 1

Primera, el servidor central elige el tipo de modelo a entrenar. En principio podría ser cualquier tipo de modelo basado en la optimización de parámetros, en nuestro caso (¡aviso de spoiler!) se usará una red neuronal.

## Bloque 2

Segunda, el servidor central transmite el modelo al resto de participantes.

## Bloque 3

Tercera, los nodos entrenan el modelo de forma local con sus propios datos.

## Bloque 4

Y cuarta, el servidor central solicita los modelos locales y a partir de ellos genera otro modelo si acceder a ningún dato.

Un ciclo completo se denomina ronda. Todo este proceso se repetirá hasta que se cumpla la condición de parada establecida, que puede estar basada en un criterio de calidad o en un número máximo de iteraciones.

## Diapositiva 9 (Diseño del experimento I)

---

Para mostrar que la aplicación del Federated Learning es una solución que puede ofrecer unos buenos resultados compararemos sus resultados con dos alternativas:

- La primera sería que cada factoría pudiese construir su propio modelo local. Sabemos que esta alternativa no es viable pero nos dará una cota superior de la calidad de los modelos que es posible construir.
- La segunda, basada en la intercambiabilidad de los modelos. Es decir, que un modelo entrenado en una planta pueda ser utilizado en otra sin falta de ser reentrenado. Notar que esta alternativa sí que cumpliría con los requisitos establecidos de velocidad de implantación y privacidad.

Para poner a prueba las tres alternativas ...

## Bloque 1

Se generarán cuatro conjuntos de datos simulados con diferentes condiciones ambientales y de funcionamiento. Estos datos contendrán tanto, información de los sensores instalados en las máquinas como de los datos de los informes de mantenimiento.

## Bloque 2

De esos conjuntos se elegirá uno (que denominaremos Piloto) y con estos datos se construirá un modelo. Este modelo será el modelo base común para todos los demás conjuntos de datos. Notar, que solo nos interesa la estructura del modelo y no el valor de sus parámetros.

## Bloque 3

Se entrenará un modelo (con la estructura definida en el paso anterior) por cada uno de los conjuntos de datos y se evaluarán los modelos obtenidos. Esto nos dará unos la base de comparación ya que nos proporcionará una cota máxima del rendimiento que podríamos esperar.

Hasta aquí la parte que concierne al entrenamiento de modelos de modo local.

## Diapositiva 10 (Diseño del experimento II)

---

En cuanto a la aplicación del aprendizaje federado ...

### Bloque 1

Se entrenará un modelo federado con tres de los cuatro conjuntos de datos. En esos tres, estará el conjunto de datos llamado Piloto. Los otros dos datasets pasarán a denominarse A y B.

### Bloque 2

Se compararán los resultados del modelo federado global con cada uno de los modelos locales de los conjuntos de datos que han participado en la federación. Es decir con los conjuntos Piloto, A y B.

### Bloque 3

Se evaluará el rendimiento del modelo federado con respecto al modelo local del conjunto de datos excluido de la federación (este conjunto se pasará a denominar N (letra inicial de New ))

Una vez examinadas todos los resultados nos será posible evaluar la idoneidad del uso del aprendizaje federado frente a su alternativa, la intercambiabilidad de los modelos.

## Diapositiva 11 (Tecnologías)

---

Para la implementación de todas las herramientas y modelos desarrollados en este trabajo ha sido necesario el uso de multitud de tecnologías. Todas ellas tienen como nexo el lenguaje de programación Python ...

### Bloque 1

¿Por qué usar Python? Se ha elegido este lenguaje por diversos motivos:

- Simplicidad. Python ofrece la posibilidad de desarrollar programas muy potentes con muy pocas líneas de código. En general, resulta un lenguaje fácil de usar y no se requiere mucho tiempo de codificación.
- Compatibilidad. Muchas de las tecnologías actuales relacionadas con el Machine Learning están pensadas para ser utilizadas con este lenguaje.
- Facilidad de aprendizaje. En comparación con otros lenguajes, Python es fácil de aprender incluso para los programadores con menos experiencia.

### Bloque 2

Todos los scripts han sido desarrollados en el entorno de desarrollo integrado (IDE) "PyCharm". Este entorno, específicamente diseñado para Python ofrece herramientas muy útiles como por ejemplo, análisis de código fuente y control de versiones.

### Bloque 3

Pandas, Pandas es una librería de para el lenguaje Python que permite la manipulación de datos y su análisis de una manera sencilla y eficiente.

### Bloque 4

Matplotlib, esta potente librería también para el lenguaje de programación Python permite realizar visualizaciones de gran calidad de una forma muy sencilla.

## Bloque 5

scikit-learn, es una librería especializada en Machine Learning para Python y contiene multitud de implementaciones de diferentes tipos de modelos. Aunque no ha sido la librería que se ha utilizado para la creación de los modelos, se ha utilizado para la evaluación de los modelos que han sido construidos con ...

## Bloque 6

... PyTorch, librería desarrollada principalmente por Facebook, de código abierto y que permite desarrollar modelos de redes neuronales profundas de una manera rápida y eficiente.

## Bloque 7

Y por último pero no menos importante PySyft. PySyft es una biblioteca de Python para el aprendizaje profundo seguro y privado. PySyft permite desacoplar los datos privados del entrenamiento del modelo, utilizando el aprendizaje federado. Esta librería es compatible con múltiples frameworks entre los que se incluyen TensorFlow y PyTorch.

## Diapositiva 12 (Simulación de un entorno industrial)

---

Para el entrenamiento de los modelos necesitamos datos. Sin embargo, es prácticamente imposible obtener un conjunto de datos real por ser este tipo de información muy sensible para las compañías. Por este motivo no se utilizarán datos reales y en su lugar, se ha desarrollado un software que nos permitirá la simulación de estos datos.

## Bloque 1

Por cuestiones obvias no nos ha sido posible construir un software que simule cada una de las diferentes máquinas que pudieran existir en un entorno industrial. Por lo tanto, en este trabajo, vamos a considerar únicamente un tipo de máquina: máquinas rotatorias de uso genérico.

## Bloque 2

Basándonos en conjuntos de datos de estudios previos para cada máquina se generarán las siguientes variables operacionales:

- Velocidad rotacional
- Temperatura
- Presión

Además, también se almacenarán datos sobre las condiciones climáticas. En nuestro caso presión y temperatura atmosférica.

## Bloque 3

También se simulará el desgaste propio del uso de las máquinas. Esto significa que, una vez lleguen al final de su vida útil, estas se estropearán. Cuando una máquina se estropee se almacenará un registro en un diario de mantenimiento tanto de cuando se produjo la avería como de cuando se realizó la reparación. Nuestras máquinas podrán fallar por dos motivos y el modelo se construirá para ser capaz de detectar el tipo de fallo.

## Diapositiva 13 (Medidas de calidad para los modelos)

---

Antes de explicar cómo se ha construido el modelo es importante hacer un pequeño inciso para hablar sobre las medidas de calidad de los modelos.

Es crucial seleccionar la medida de evaluación adecuada para cada tipo de problema. Por ejemplo, un clásico error es considerar la medida "accuracy" (número de casos correctamente clasificados entre número total de casos) para problemas con conjuntos de datos desbalanceados donde, un número cercano al 1 (valor máximo) no indica en ningún caso que el modelo sea bueno ya que, un modelo cuya salida sea siempre la clase mayoritaria. Este modelo ofrecerá valores altos "accuracy" lo cual no significará en ningún caso que este sea un buen modelo.

Aunque existen muchas más en esta presentación hablaremos únicamente de las medidas que usaremos en este trabajo.

## Bloque 1

Precision, esta medida responde a la siguiente pregunta ¿Qué proporción de instancias clasificadas como "X" son realmente "X"? Se busca maximizar esta medida cuando queremos estar muy seguros de nuestra predicción.

En nuestro caso la maximización de esta medida nos ayudaría a estar seguros de que un error se va a producir con un alto grado de certeza, esto ahorraría dinero a la empresa ya que solo se harían mantenimientos cuando fueran estrictamente necesarios, el problema con esto es que podríamos dejar pasar fallos de los que no estamos muy seguros y que potencialmente representarían un peligro para la seguridad de la instalación.

## Bloque 2

Recall, esta medida responde a la siguiente pregunta ¿De todas las instancias "X" que existen, qué proporción han sido clasificadas como "X"? Se busca maximizar esta medida cuando queremos capturar la mayor cantidad de clases "X" posible.

La maximización de esta medida es equivalente a la maximización de la seguridad en la instalación ya que nos garantiza que todos los casos en los que se produzca un fallo serán detectados. Esto tiene su contrapartida ya que esto puede ocasionar falsos positivos lo que desembocaría en un aumento en los costes de mantenimiento.

## Bloque 3

Aunque obviamente la seguridad siempre debe ser lo primero, es necesario establecer un compromiso entre las dos medidas. Para evaluar este compromiso se utiliza la medida f1-score, que consiste en la media armónica de las dos medidas anteriores.

## Diapositiva 14 (Construcción del Modelo Base)

---

Después de este pequeño inciso pasemos a describir la construcción del modelo base.

Recordemos que el objetivo es poder decidir si en un determinado momento una máquina está en riesgo de rotura o no utilizando únicamente sus datos telemétricos. Dicho de otra manera, se debe clasificar esa máquina como potencialmente peligrosa o como segura. Teniendo en cuenta esto parece claro que será necesario el uso de un modelo de clasificación.

¿Cómo hemos construido el modelo?

## Bloque 1

Primero se ha procedido a la preparación de los datos para la tarea de clasificación. Esta preparación la podemos dividir en 4 fases:

### Agregación

Muchas máquinas del mundo real funcionan en ciclos. Un ciclo puede considerarse como un período temporal que describe un estado de funcionamiento de una máquina. Por ejemplo, la operación de un motor en un avión puede describirse por los ciclos: motor en funcionamiento (avión en vuelo) o motor apagado (avión en tierra).

Las transmisiones de telemetría sin procesar, si bien pueden ser muy útiles para tareas como el monitoreo en tiempo real, pueden causar problemas a la hora de construir modelos de detección de fallos. Es frecuente que en entornos no controlados (como

puede ser una fábrica) los datos no sean del todo precisos. Para mitigar los posibles errores en las mediciones (fallos puntuales en los equipos de medida, ruido, etc.) que pueden afectar a la calidad de los modelos construidos, suele ser una buena opción utilizar datos agregados por ciclo. Las funciones de agregación que se han sido utilizadas han sido la media y el máximo.

### Etiquetado

Gracias a los datos de mantenimiento podemos saber cuándo una máquina ha fallado, por lo tanto, podemos etiquetar el registro telemétrico en ese instante como un fallo (la etiqueta dependerá del tipo de fallo recordemos que podían ocurrir dos tipos de fallo). Como es interesante detectar el fallo antes de que se produzca, para tareas de mantenimiento preventivo, por ejemplo, se etiquetaran también como fallos los N registros anteriores al fallo real (nosotros hemos usado  $N = 7$ )

### Enriquecimiento

Una forma de añadir mayor cantidad de información al conjunto de datos disponible es lo que se conoce como ingeniería de características. La ingeniería de características consiste en la creación de nuevos atributos a partir de los ya existentes. En muchos casos este tipo de procedimientos mejora notablemente la calidad de los modelos obtenidos.

Aquí se añadirán nuevos atributos basados en los atributos ya existentes. Cada nuevo atributo se corresponderá con los datos promediados de su atributo asociado en las N instancias precedentes (nosotros usaremos  $N = 5$ ).

Este procedimiento es muy interesante ya que permite añadir una cierta cantidad información histórica a cada registro. El modelo podrá, no solo tener información instantánea si no que tendrá también información sobre la tendencia.

### Balanceo

Por la propia naturaleza de los datos que estamos manejando es normal que exista una gran diferencia en número entre los casos donde no se detecta ninguna anomalía (clase 0) y los casos donde es posible que se produzca una avería (clase 1 o clase 2). Sin embargo, este tipo de conjuntos de datos suelen ser problemáticos a la hora de entrenar los modelos. Existen multitud de técnicas para balancear conjuntos de datos. Nosotros utilizaremos el algoritmo SMOTE. Este algoritmo permite generar nuevas instancias a partir de las clases minoritarias dejando intacta la clase mayoritaria. Las nuevas instancias no son copias de los casos existentes si no que se calculan como combinaciones lineales de los vecinos más cercanos de esa misma clase.

## Bloque 2

Elección del modelo. Hay muchas de familias de modelos que pueden ser usados para tareas de clasificación: árboles de decisión, máquinas de soporte de vectores, k-vecinos más cercanos, etc. Sin embargo nos hemos decidido por el uso de redes neuronales. Esta elección está motivada principalmente por dos cuestiones:

1. Las redes neuronales han demostrado tener un rendimiento excelente en multitud de problemas.
2. Aunque en teoría el uso del Federated Learning es aplicable a cualquier tipo de modelo cuyo entrenamiento se base en la optimización de parámetros, es cierto que al ser una tecnología relativamente nueva la mayoría de los frameworks actuales solo permiten el uso de redes neurales en sus implementaciones de aprendizaje federado.

## Bloque 3

Una vez elegido el tipo de modelo será necesario determinar su fisonomía (número de capas, neuronas por capa, etc.) entrenarlo y validarlo. Todas estas operaciones se han llevado a cabo siempre de acuerdo a las medidas de calidad mostradas en la anterior transparencia.

Es importante destacar que:

1. El modelo ha sido construido únicamente con el conjunto de datos denominado "Piloto"
2. Del modelo aquí construido solo no interesa su estructura no sus parámetros ya que el proceso de entrenamiento se llevará a cabo para cada conjunto de datos de forma individual.
3. De aquí en adelante todos los modelos tendrán la misma estructura y serán entrenados de la misma forma.



## Diapositiva 15 (Resultados del modelo base)

---

Una vez construido el modelo base, evaluamos su rendimiento en los distintos conjuntos de datos de los que disponíamos. Estos valores nos servirán como referencia a la hora de comparar las diferentes alternativas que se han propuesto como solución al problema.

Se puede observar que todas las plantas arrojan valores similares en cuanto a: precision, recall y f1-score, en torno a un 90, 95 y 94 por ciento respectivamente. En las plantas "Piloto" y "A" se observa una diferencia bastante significativa entre el recall y la precision mientras que en las otras dos: la Planta B y la Planta N los valores están más cercanos. Cabe destacar que los resultados obtenidos en la Planta N son los peores de entre los cuatro datasets.

## Diapositiva 16 (Intercambiabilidad del modelo base)

---

En la sección anterior se ha visto lo que se puede esperar del modelo base si es entrenado de manera local. En esta diapositiva se pretende responder a la siguiente pregunta: ¿Funciona bien un modelo entrenado en una planta en otra sin necesidad de reentrenarlo? Para resolver esta duda presentamos a continuación una tabla comparativa del f1-score (media de todas las clases) sobre todas las posibles combinaciones train-test. La razón de elegir esta medida es por simplicidad, ya que con un único valor podemos hacernos una idea tanto de la precision como del recall.

Observando la tabla podemos ver que por norma general el modelo entrenado de forma específica (diagonal) supera ampliamente en rendimiento a los entrenados en otras instalaciones. Es destacable el caso de la instalación "N" en la que para algunos casos el score es mucho mejor que para su propio conjunto de entrenamiento. El motivo seguramente guarde relación con el hecho de que el número de rondas de entrenamiento no fuera lo suficientemente alto. Aunque este hecho aislado de tener puntuaciones altas con modelos ajenos, fuera común, existiría otro problema, ¿Cómo saber a priori que modelo de todos de los que se dispone es mejor para la planta objetivo?

Teniendo en cuenta los datos y reflexiones anteriores no parece posible muy viable traspasar un modelo de una instalación a otra sin que esto tenga una repercusión negativa en la calidad de las predicciones del modelo.

## Diapositiva 17 (Resultados del modelo federado)

---

No entraremos en detalles de la construcción del modelo federado por ser esta, análoga a la construcción de los modelos locales. Si embargo, si que queremos destacar que las instalaciones que han participado en la federación (indicadas con un icono a su izquierda) han sido las plantas: "Piloto", "A" y "B". Analicemos ahora los resultados obtenidos.

Vemos que, aunque las medidas de calidad de no son tan buenas como los del entrenamiento de forma local vemos que se mantienen en unos valores aceptables. Nuevamente debemos destacar los valores, bastante bajos de la precisión en el caso de la instalación N. Esto seguramente se deba a que:

1. En el modelo base de esa instalación ya presentaba valores más bajos que en resto de las instalaciones.
2. Que el modelo no ha participado en la federación a la hora de ser entrenado.

## Diapositiva 18 (Comparación entre aproximaciones)

---

Veamos ahora una comparación entre todos los modelos. Al igual que anteriormente se utilizará únicamente la medida f1-score por simplicidad.

Como es lógico ambas alternativas presentan resultados peores al entrenamiento local que sería el caso, digamos, ideal. Por otro lado, se puede ver que el método de aprendizaje federado es siempre mejor que el peor de los casos cuando se utiliza otra instalación, incluso en algunos casos supera a la mejor de las opciones. Hay que tener en cuenta que, aunque en el mejor de los casos de usar el modelo de otra instalación supera al aprendizaje federado, nos encontramos con el problema adicional de encontrar, a priori, cual de todas las instalaciones disponibles será la más adecuada. Por lo tanto, parece claro que el uso del Federated Learning es muy deseable en casos de uso como los descritos en este trabajo.

## Diapositiva 19 (Conclusiones)

---

Recapitulando, el objetivo de proyecto era explorar el posible uso del Federated Learning para la detección de eventos anómalos dentro de un entorno industrial. Para ello se ha descrito un escenario que podría corresponderse con las necesidades de una compañía multinacional como podría ser una compañía siderúrgica, minera, un fabricante de productos químicos, etc. Se han expuesto las limitaciones existentes en cuanto a la distribución de datos entre instalaciones en el ámbito de una organización con una gran dispersión geográfica y se ha puesto de manifiesto la necesidad de una rápida implantación de modelos de Machine Learning en instalaciones de nueva creación.

Para la solucionar el problema presentado se han comparado dos soluciones:

- Una basada en la intercambiabilidad de modelos
- Una basada en el uso del Federated Learning

### Bloque 1

Se ha mostrado que el método basado en la intercambiabilidad entre modelos y el basado en aprendizaje federado ofrecen resultados similares (considerando siempre el mejor de los casos) pero añade complejidad al problema. Es necesario crear un método para decidir la planta de origen del modelo ya que como se ha visto, no todos los modelos ofrecen los mismos resultados.

### Bloque 2

En cualquier caso el modelo federado ofrece siempre mejores resultados que la peor de las soluciones de intercambiabilidad.

### Bloque 3

Otro problema que habría que considerar es la propiedad del modelo, una planta podría exigir a otra algún tipo de contrapartida por la cesión del modelo creado con sus datos. En el caso del Federated Learning todos los participantes son responsables en la creación del modelo por lo que nadie es propietario exclusivo de este.

Teniendo en cuenta todos estos factores, se ha demostrado que el aprendizaje federado es una excelente forma de construir modelos que puede ser aplicada en muchos ámbitos industriales en los que, el acceso a datos sea restringido y en donde sea necesaria una rápida implantación y que además ofrece mejores resultados que otras alternativas.

## Diapositiva 20 (Trabajo futuro)

---

El objetivo del trabajo es mostrar una metodología por lo que el alcance de los modelos aquí generados no se extiende más allá de su uso a nivel didáctico y su aplicación en entornos reales dependerá mucho del tipo de entorno y de las fuentes de datos disponibles. Si embargo, todo el procedimiento hasta llegar a su construcción puede resultar de gran interés en la resolución de problemas similares y justamente esto, el proceso, es lo que deberá considerarse como el producto final de este trabajo. Posibles líneas de trabajo futuro podrían ser la aplicación de los métodos aquí descritos en entornos industriales reales con conjuntos de datos reales.

## Diapositiva 21 (Gracias)

---

Y aquí finaliza la presentación. Gracias por su atención.