

# Machine Learning para preservar la anonimidad

Normalmente las técnicas de machine learning deben manejar volúmenes de data enormes, lo que requiere muchos recursos de computación, el uso de nodos independientes no es apto debido a límites de memoria o de tiempo. Algunos de los frameworks más usados para machine learning están optimizados para usar GPUs. Existen muchos sistemas distribuidos de computación capaces de realizar entrenamiento en múltiples nodos de forma coordinada. Las aproximaciones convencionales requieren de una plataforma centralizada que recoge los datos y los distribuye entre los nodos.

## Técnicas de preservación de privacidad en Machine Learning

- **Computación Multi-Grupo Segura:** Una función puede ser computada colectivamente por varios grupos sin mostrar sus propios datos. Cada grupo debe intercambiar sus salidas con otros en secreto para que puedan ser agregados. Es una de las formas más seguras de entrenar modelos de machine learning. Por desgracia no es bueno para entrenar modelos muy complejos o a gran escala.
- **Cifrado homomorfo:** Puede ser aplicada directamente para cifrar datos que luego son transferidos a un servidor central. Cualquier modelo puede, en teoría, ser entrenado con estos datos cifrados, obteniéndose un modelo cifrado que puede ser enviado a los clientes. El server no es capaz de descifrar ni los datos ni el modelo. Los clientes tienen acceso a sus propios datos y al modelo descifrado una vez entrenado. Este modelo tampoco es adecuado para larga escala.

## Aprendizaje Federado

Es un procedimiento distribuido y colaborativo para entrenar modelos de Machine Learning sin mostrar los datos con los que se entrena. La idea de este método es mover la computación al borde, es decir, el dispositivo que obtiene los datos. Los datos no son nunca transferidos a un servidor central o a ningún almacenamiento centralizado, manteniéndose aislados localmente. Esto mitiga vulnerabilidades relacionadas con los datos. El machine learning es un modelo non-IID (Non-Independent identically distributed data).

Los algoritmos de aprendizaje están descentralizados, confiando en cada nodo para realizar entrenamiento parcial del modelo. Cada nodo computa actualizaciones del modelo parcial con sus propios datos, intercambiando los parámetros con otros grupos. Un servidor central suele ser requerido para coordinar el entrenamiento, añadiendo los resultados del entrenamiento para calcular el modelo global. Una vez entrenado, el modelo global se distribuye a cada nodo para que lo pueda usar para realizar predicciones o más iteraciones del modelo.

## Problemas de privacidad

Los datos no salen del dispositivo del usuario, por lo que el acceso directo por parte de terceras partes no es posible, pero sigue siendo posible la mala utilización de los datos y los modelos de Machine Learning pueden contener información sensible sobre los datos, siendo posible realizar ataques de inferencia.

### Ataques de inferencia en aprendizaje federado

- Ataques de inversión de modelo
- Ataques de inferencia de miembros
- Tipos de objetivo:
  - Ataque de Caja negra (pasivo): El objetivo es el modelo final ya entrenado
  - Ataque de Caja blanca (activo): Monitoriza los cambios del modelo en cada ronda del entrenamiento
- Tipo de atacante:
  - Cliente: Puede inspeccionar las versiones consecutivas del modelo global sin interferir con el procedimiento
  - Coordinación del lado del servidor: Inspecciona las actualizaciones parciales del modelo enviada por los clientes.

### Ataques de envenenamiento en Aprendizaje Federado

From:  
<https://www.knoppia.net/> - Knoppia

Permanent link:  
[https://www.knoppia.net/doku.php?id=pan:machine\\_learning\\_privacy\\_v2&rev=1767821918](https://www.knoppia.net/doku.php?id=pan:machine_learning_privacy_v2&rev=1767821918)

Last update: **2026/01/07 21:38**

