

UNIVERSIDAD
COMPLUTENSE
DE MADRID



Autor: María José
Gómez Silva

Machine Learning con Python. Semana 1.

CAPÍTULO 1. INTRODUCCIÓN AL MACHINE LEARNING.

1.7 Etapas de un proyecto de Machine Learning

La primera etapa cuando se trabaja con aprendizaje automático es la de *preparación de los datos*, también llamada *extracción de características*, como puede observarse en el diagrama de bloques de la Figura 1. En esta etapa los datos pueden ser filtrados, normalizados, escalados, seleccionados y en general, procesados. Esta etapa trata de mejorar los datos o seleccionar aquellos más representativos para la realización de la tarea en cuestión. En esta etapa se realiza el procesamiento necesario para adecuar los datos de modo que el proceso de entrenamiento sea favorecido, facilitado o mejorado. Algunos de los métodos más usuales en esa primera etapa serán abordados en el siguiente capítulo de *Ingesta de Datos*.

Cuando se realiza un aprendizaje supervisado, además, es necesaria una etapa de *etiquetado* o *anotación*, para determinar manualmente (o con algún mecanismo semiautomático) la solución o etiqueta de cada muestra empleada como dato de entrada en el entrenamiento. Es decir, para los datos de entrenamiento, es necesario anotar la predicción que esperamos que el modelo dé cuando haya aprendido. No sólo los datos de entrenamiento deben ser anotados, también los de validación y test, como se verá a continuación.

Una vez preparados los datos, el proceso de *entrenamiento* o *aprendizaje* automático es un proceso cíclico en el que el modelo es actualizado en función de la desviación de sus predicciones con respecto a la solución deseada, como se puede apreciar en la figura.

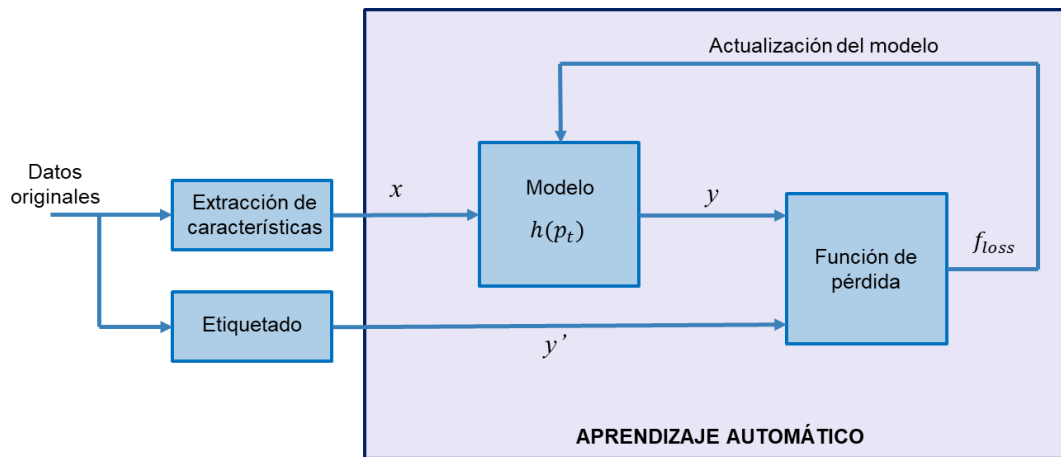


Figura 1. Esquema del proceso de entrenamiento

Por lo tanto, aprender es actualizar los parámetros del modelo de modo que en cada ciclo se desvíe menos de su objetivo.

Tras el proceso de aprendizaje, el modelo debe ser evaluado. La *evaluación* o *test* se realiza calculando la precisión, o cualquier otra métrica que permita medir la bondad del modelo. En este proceso, el modelo no es actualizado, como se observa en la Figura 2.

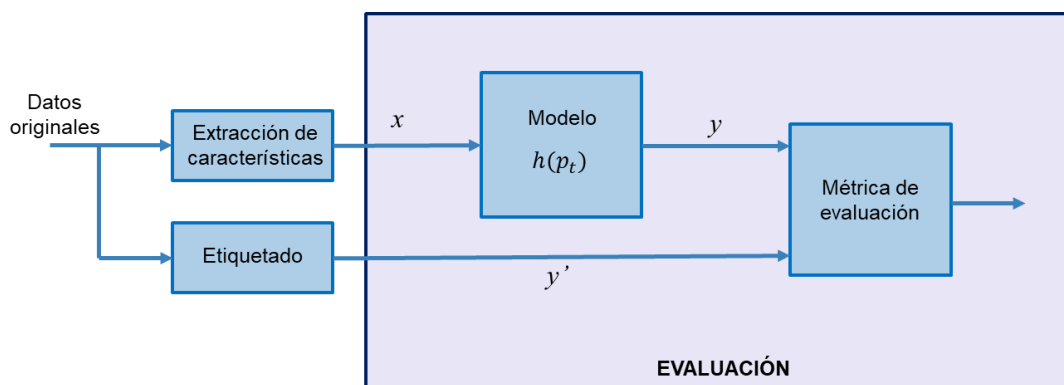


Figura 2. Esquema del proceso de evaluación

Una práctica muy común es realizar algunas evaluaciones del modelo durante su entrenamiento, cada un cierto número de iteraciones de aprendizaje. Estas evaluaciones, puntuales y periódicas durante el entrenamiento, se denominan validaciones. La validación permite observar cómo progresa la precisión del modelo para datos desconocidos para el modelo, distintos a los de entrenamiento.

En resumen, un proyecto de Machine Learning requiere tres conjuntos de datos diferentes: el de entrenamiento, el de validación y el de test. En caso de realizar un aprendizaje supervisado, todos los datos deben estar etiquetados, no sólo los datos de entrenamiento, ya que, para estimar la métrica o métricas de evaluación, es necesario comparar la predicción del modelo con respecto al valor esperado de referencia.

El rendimiento del modelo sobre los datos de entrenamiento determina el ajuste de los parámetros del modelo. El resultado sobre los datos de validación puede emplearse para ajustar algunos hiper-parámetros como el *learning rate*, o para decidir el final del entrenamiento. Finalmente, el rendimiento del modelo entrenado es evaluado sobre el conjunto de datos de test.

Una vez finalizado el aprendizaje del modelo, y habiendo obtenido una evaluación positiva del mismo, el modelo puede ser empleado para realizar la tarea para la que se le ha entrenado, pero ahora sobre nuevos datos, de los cuales no se conoce la solución a priori, sino que es el modelo quien debe estimarla, como se observa en la Figura 3.

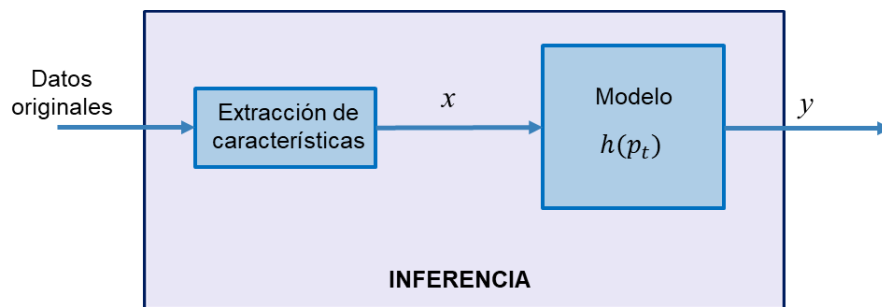


Figura 3. Esquema de inferencia con un modelo

En un entorno de experimentación real, las etapas de preparación de los datos, entrenamiento y evaluación son a su vez iterativas. Por ejemplo, de la evaluación de distintos modelos entrenados con distintas características de los datos se pueden obtener conclusiones sobre qué características mejoran o no la tarea. Por ello, el aprendizaje automático es un proceso interdisciplinar, cuyo éxito depende de la colaboración de expertos en distintas áreas: expertos en la tarea que se esté intentando modelar, en adquisición y análisis de datos, en algoritmos de “Machine Learning”, y en implementación de los modelos finales, entre otras.