

Conceptos básicos de las máquinas de aprendizaje

Conceptos básicos de las máquinas de aprendizaje	1
Competencias	2
Introducción	2
Précis: Conceptos básicos de Machine Learning (Aprendizaje de Máquinas)	3
Regresión	3
Clasificación	3
El proceso de entrenamiento	4
Subajuste y sobreajuste	5



¡Comencemos!

Competencias

- Conocer la terminología y conceptos básicos ocupados en machine learning.
- Tener un primer acercamiento a los problemas canónicos en máquinas de aprendizaje.

Introducción

En esta sección conocerás la terminología y conceptos básicos ocupados en machine learning, además de tener un primer acercamiento a los problemas canónicos en máquinas de aprendizaje

¡Vamos con todo!



Précis: Conceptos básicos de Machine Learning (Aprendizaje de Máquinas)

Antes de comenzar a ver modelos y funciones matemáticas, es bueno dedicar tiempo a asimilar los conceptos básicos que serán utilizados en este módulo y que encontrarán en las referencias y artículos sobre Machine Learning en general.

En Machine Learning hay dos problemas canónicos dentro de los cuales se pueden clasificar casi la totalidad de los modelos actuales: El problema de la regresión y el de la clasificación.

Regresión

Hace referencia a la habilidad de un modelo de ser capaz de entregarnos un output en forma numérica. Por ejemplo, predecir el valor de la temperatura para el día siguiente, predecir el valor de una acción en la bolsa o incluso la probabilidad de algún suceso. Todo método que entregue un output numérico se considera en general que está resolviendo un problema de regresión.

Clasificación

Son problemas donde se quiere encontrar o predecir la clase a la que pertenece un cierto registro. Por ejemplo, dada la descripción física de una persona (peso, estatura, medidas corporales) determinar si es hombre o mujer, o determinar la acción más conveniente a ser realizada frente a ciertas condiciones. En todos estos casos, el output no es un número, sino una clase.

Con lo anterior en mente, podemos definir una máquina de aprendizaje (o un learner) como un modelo matemático con los siguientes componentes:

- **Una Función Objetivo:** Es la función matemática que el modelo tratará de optimizar, al optimizar esta función el modelo se entrena y aprende los parámetros correspondientes.
- **Entrenamiento:** Es el proceso mediante el cual la máquina optimiza su función objetivo y, al hacerlo, encuentra los valores de sus parámetros. Cada máquina distinta tiene parámetros distintos, por lo tanto para un mismo problema, resolverlo con dos o más máquinas distintas significa entrenar ambas máquinas por separado.
- **Espacio de parámetros:** Es el espacio en el cual se encuentran todos los posibles parámetros para la máquina que estamos implementando. El espacio de atributos no es lo mismo que el espacio de parámetros.

- **Atributos y variable objetivo:** Un atributo es una dimensión de medición, una columna en nuestro dataset. Por ejemplo, la estatura de una persona es un atributo, la edad es otro atributo dentro de una base de datos. La variable objetivo (o target) es la variable que queremos predecir, ésta puede ser numérica (regresión) o categórica (clasificación), por ejemplo, el sexo de una persona (categórica) o la edad de la misma (numérica).

El proceso de entrenamiento

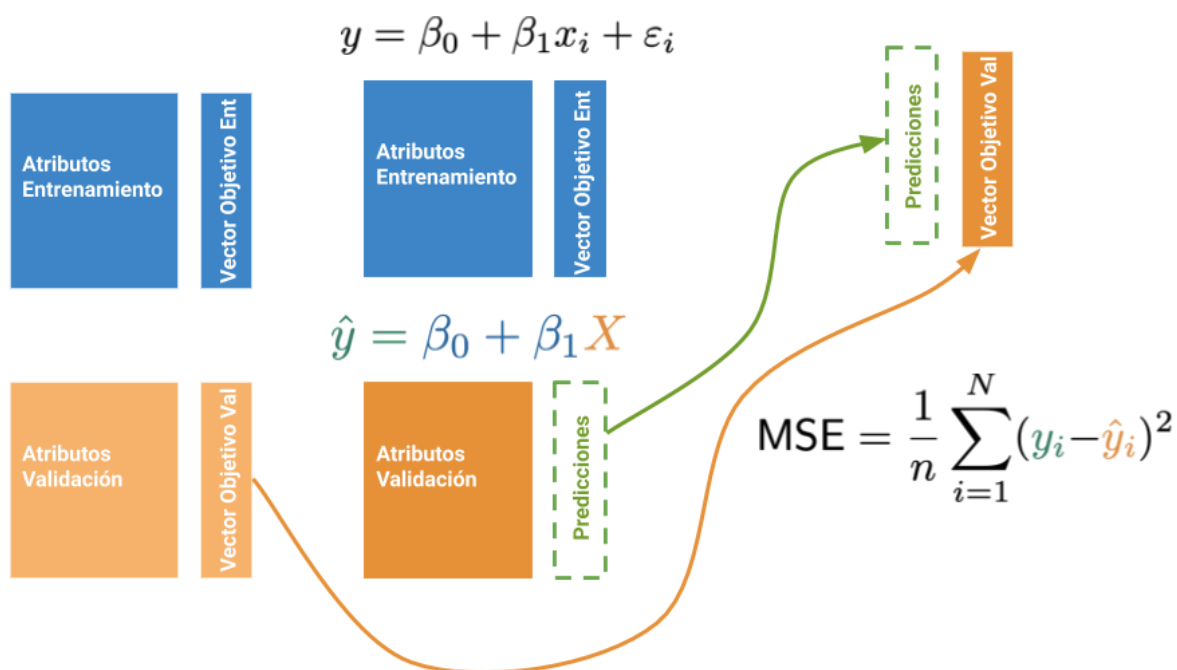


Imagen 1. Proceso de entrenamiento.
Fuente: Desafío Latam.

Una vez definido los elementos básicos para implementar una máquina de aprendizaje, definiremos los pasos estándares para obtener un learner efectivo.

Se puede resumir en los siguientes pasos:

- **División de la muestra:** Mediante la división de la muestra en conjuntos de entrenamiento y validación nos aseguramos de poder "replicar" el comportamiento de nuestro modelo en un nuevo conjunto de datos.
- **Entrenamiento del modelo:** De esta manera, generaremos el entrenamiento de un modelo candidato que tendrá una combinación de parámetros que permitan generar predicciones.

- Evaluación del modelo: Para evitar vicios en el entrenamiento, parcelamos un conjunto de datos para el comportamiento del modelo en un nuevo conjunto de datos.

Subajuste y sobreajuste

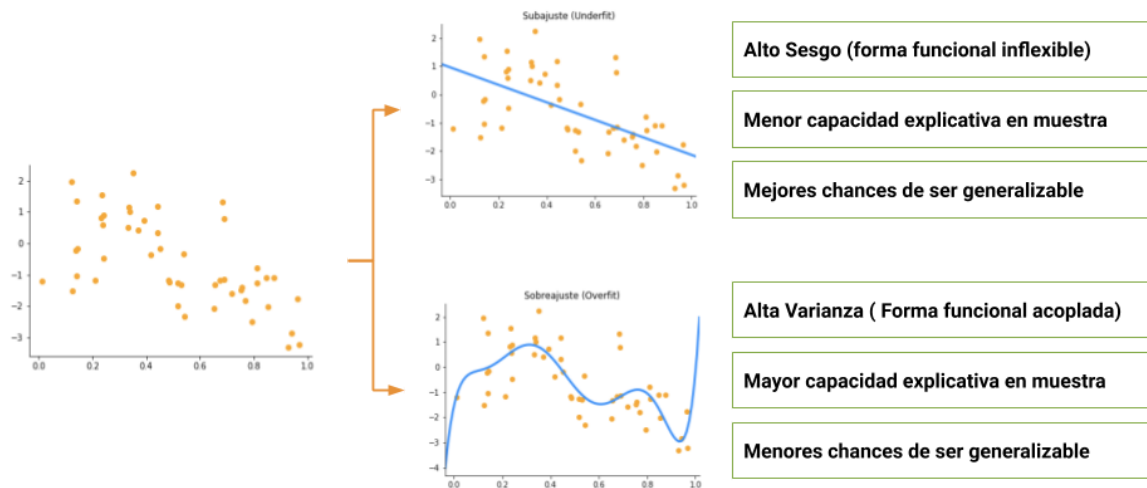


Imagen 2. Subajuste y sobreajuste.

Fuente: Desafío Latam.

- **Overfitting (sobre-ajuste):** Fenómeno en el cual nuestro modelo aprendió demasiado bien el conjunto de datos de entrenamiento. Como consecuencia de esto, el modelo generó reglas internas que se apegan demasiado a estos datos y al evaluarlo en datos que nunca antes ha visto (datos nuevos), se comporta mal. Un ejemplo simple para recordar esta definición es pensar en un mal estudiante que aprende de memoria como hacer los ejercicios de la guía, sin entender en realidad el fenómeno que ocurre por detrás, como es de esperarse, cuando llega el momento de la prueba/certamen, no es capaz de generalizar lo poco que sabe y le va mal. Usualmente evitamos el overfitting implementando técnicas de regularización.
- **Underfitting (sub-ajuste):** Fenómeno en el cual nuestro modelo aprendió muy poco sobre el fenómeno subyacente en los datos y no es capaz de generalizar lo que aprendió durante el entrenamiento a datos nuevos. Usualmente evitamos el underfitting implementando técnicas de data augmentation o simplemente recolectando más data.