



Programa Profesional en Inteligencia Artificial y Data
Science

Módulo 3. Aprendizaje Automático Aplicado

Tema 1. Introducción al lenguaje automático

1.1. Introducción y objetivos

1.2. Introducción

Tema 2. Análisis exploratorio de datos y preprocesamiento

2.1. Introducción y objetivos

2.2. Análisis exploratorio de datos y preprocesamiento I

2.3. Análisis exploratorio de datos y preprocesamiento II

Tema 3. Algoritmos de aprendizaje automático supervisado

3.1. Introducción y objetivos

3.2. Algoritmos de aprendizaje automático supervisado I

3.3. Algoritmos de aprendizaje automático supervisado II

Tema 4. Algoritmos de aprendizaje automático no supervisado

4.1. Algoritmos de aprendizaje automático no supervisado I

4.2. Algoritmos de aprendizaje automático no supervisado II

4.3. Aprendizaje automático

Tema 5. Redes neuronales y deep learning

5.1. Introducción y objetivos

5.2. Redes neuronales y Deep learning (continuación)

Tema 6. Despliegue de modelos de aprendizaje automático

6.1. Introducción y objetivos

6.2. Despliegue de modelos de aprendizaje automático

Tema 1. Introducción al lenguaje automático

1.1. Introducción y objetivos

En esta unidad formativa se hace una introducción al aprendizaje automático. Inicialmente describimos y caracterizamos la inteligencia artificial fuerte y débil.

Luego se introducen los conceptos principales de aprendizaje automático y sus tipos para, a continuación, ir esbozando los conceptos de conjunto de datos, los tipos de variables y su clasificación para terminar haciendo un recorrido por los algoritmos y modelos de aprendizaje automático y su proceso.

Tema 1. Introducción al lenguaje automático

1.2. Introducción

Probablemente ya conozcas de lecturas anteriores, por videos u otra vía, los elementos esenciales de la inteligencia artificial. Probablemente sepas también, que se tiene como inicios de este campo de las ciencias de la computación, el año 1956 cuando se juntaron, en lo que se dio a llamar conferencia de Dartmouth en Estados Unidos, científicos de la talla de John McCarthy, Marvin Minsky, Claude Shannon, Allen Newell, y Herbert. A. Simon. En esta conferencia se acuñó el término «inteligencia artificial» para denominar el nuevo campo que acaba de nacer.

El objetivo de dicha conferencia era examinar la posibilidad de que cada aspecto del aprendizaje o cualquier otra característica de la inteligencia pudiera ser descrita con tanta precisión que se pudiera construir una máquina que la simulara. Se intentaría encontrar la forma para que las máquinas usen el lenguaje, creen abstracciones y conceptos, resuelvan los tipos de problemas que actualmente solo pueden solucionar los humanos y se vayan autoperfeccionando. Pensaron que diez personas tardarían dos meses en conseguirlo. Sin embargo, conseguir ese objetivo no ha sido tan sencillo como inicialmente pensaban estos pioneros de la IA. La imposibilidad de cumplir este objetivo llevó al filósofo John Searle a concebir la IA como dividida en dos bandos, dos partes, a los cuales los llamó IA fuerte y IA débil.

Para él, la IA fuerte es la que se dedica, aún, a perseguir el objetivo inicial de la IA. Aquel que pretende desarrollar máquinas capaces de comportarse como un humano más, este campo también se ha nombrado Inteligencia Artificial General.

Por otra parte, la inteligencia artificial débil, es aquella que se centra en resolver problemas más concretos o estrechos, tales como el reconocimiento del habla, la traducción, el reconocimiento de imágenes, etc.

Tema 1. Introducción al lenguaje automático

En la siguiente tabla mostramos una caracterización de cada una de estas Inteligencias Artificiales.

Inteligencia artificial débil	Inteligencia artificial fuerte
Existe en la vida Real	Solo existe en la ciencia ficción
AlphaGo, Watson, etc	3CP0, Wall-e, etc.
Orientada a Problemas muy concretos. Por ejemplo, aprender a clasificar imágenes.	Orientada a resolver problemas abiertos. Ejemplo, modelar la conciencia o el sentido común.
Reactivo	Proactivo
Inflexibles. Solo funcionan para el problema que se diseña	Flexibles. Se adaptan a cualquier situación. Aprenden de sus límites y posibilidades
Son programadas o entrenadas por humanos	Se autoprograman
Implementadas sobre miles de neuronas artificiales. (Deep Learning por ejemplo que pueden tener cientos de capas ocultas)	Implementadas sobre miles de millones de redes o neuronas artificiales u otras tecnologías. (Computación cuántica)
No razonan, solo computan	Imitan el comportamiento humano
Aprenden a partir de ejemplos similares	Aprenden como los humanos
Tareas repetitivas	Aprenden nuevas tareas
No pueden salirse de su marco de trabajo	Se adaptan a nuevos entornos de trabajo

Tabla 1. Caracterización de IA. Fuente: elaboración propia.

Uno de los nombres que representa de forma inequívoca la persecución de los objetivos iniciales de la IA y por tanto representa a los investigadores de la IA fuerte es el del investigador Ben Goertzel.

En un intento por desarrollar lo que se conoce como IA fuerte, Goertzel quería crear un cerebro de bebé digital y lanzarlo a Internet, donde creía que crecería para volverse completamente consciente de sí mismo y mucho más inteligente que los seres humanos. «Estamos en vísperas de una transición de igual magnitud a la aparición de la inteligencia o al surgimiento del lenguaje», dijo al Christian Science Monitor en 1998.

Tema 1. Introducción al lenguaje automático

A pesar de los esfuerzos, la IA Fuerte aún no se ha logrado y no se tiene claro cuando se logrará. Al punto que el padre de la Inteligencia Artificial John McCarthy, casi al final de su carrera investigadora declara «Para crear una verdadera IA se necesitaría el trabajo de 1,7 Einsteins, 2 Maxwells, 5 Faradays y la financiación de 0,3 Proyectos Manhattan», en una clara alusión a la imposibilidad de conseguirla.

Sin embargo, mientras la IAG, conseguía poco éxito, la IA débil alcanzaba hitos importantes. A continuación, se expone una lista, que, sin ser exhaustiva, muestra verdaderos logros de este campo.

Algunos hitos de la IA débil:

- ▶ 1979 el programa BKG 9,8 gana al campeón mundial de Backgammon.
- ▶ 1980 aparece la primera furgoneta guiada por visión artificial.
- ▶ 1988 se aplica la IA a la traducción de Inglés a francés.
- ▶ 1994 dos vehículos autónomos recorren 1000 km en las autopistas de París.
- ▶ 1996 nacen los agentes inteligentes que perciben el entorno.
- ▶ 1997 la IA Deep Blue gana al campeón mundial de Ajedrez y gran maestro Garry Kasparov.
- ▶ 2004 vehículos autónomos recorren 100 kilómetros en el desierto de Mojave en una competencia.
- ▶ 2008 Google lanza la primera aplicación de reconocimiento de Voz.
- ▶ 2011 Watson, un programa desarrollado por IBM, gana el concurso Jeopardy a contrincantes humanos.
- ▶ 2012 aparecen las primeras aplicaciones de Aprendizaje computacional que son capaces de reconocer y clasificar imágenes con gatos.

Tema 1. Introducción al lenguaje automático

- ▶ 2013 Boston Dynamics construye su robot Atlas, bípedo de rescate.
- ▶ 2014 el programa Eugene Goostman pasa el test de Turing. (El 33 % de los jueces creyeron que hablaban con un humano)
- ▶ 2016 AlphaGo gana al campeón mundial de Go. El MIT lanza el primer Taxi autónomo en Singapur.
- ▶ 2017 el programa Libratus vence a cuatro oponentes jugando Poker. La experiencia adquirida por Deepmind puede ahora ser transferida a otro juego. La Ginoide Sophia, se convierte en ciudadana saudí. (Ginoide, término utilizado para denominar a un androide de aspecto antropomorfo femenino).
- ▶ 2018 Alpha Zero aprende a jugar Ajedrez por sí misma. El robot bípedo de Boston Dynamic aprende a hacer Parkour.

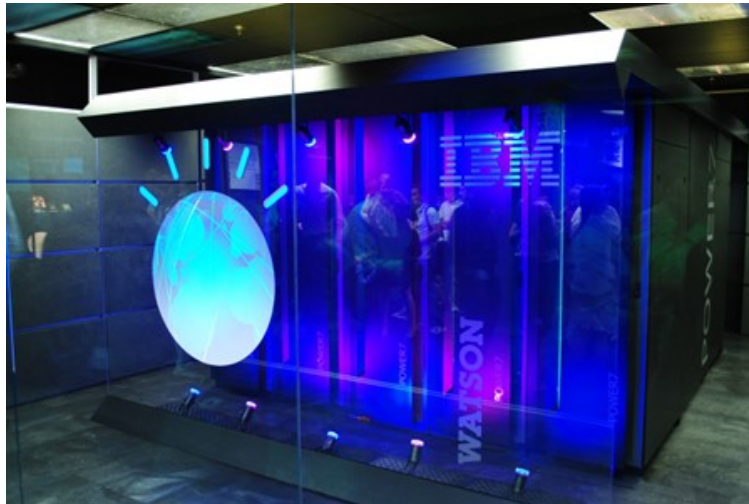


Figura 1. Soporte de Hardware del programa de IBM, Watson.

Tema 1. Introducción al lenguaje automático



Figura 2. Atlas, el robot de Boston Dynamic.

El problema del aprendizaje automático

Uno de los elementos esenciales de la inteligencia humana, es la capacidad de aprender. El interés por el aprendizaje y el razonamiento no es nuevo. Ya desde Sócrates y Aristóteles se estudia de forma analítica la inferencia deductiva e Inductiva y de manera general el aprendizaje en animales y humanos. El desarrollo de los ordenadores ha permitido contar con una nueva herramienta para su investigación y ha engendrado un nuevo campo de investigación, el aprendizaje automático.

El aprendizaje tiene que ser visto como un proceso y uno de los elementos centrales de este es el razonamiento. Este es, a su vez, un proceso del pensamiento por medio del cual, a partir de determinadas premisas, se llega a conclusiones y por tanto se genera conocimiento.

Si las premisas son casos concretos, particulares, que se comparan entre sí, y determinando lo que le es común y ese algo que es común no es casual sino esencial a ellos y por tanto podemos llegar a conclusiones tales como «Todos los

Tema 1. Introducción al lenguaje automático

elementos de tipo x tienen la cualidad y » estamos en presencia de razonamiento inductivo y de esta forma se descubren leyes y principios. La principal cualidad del razonamiento inductivo es que parte de premisas puntuales, particulares y llega a conclusiones generales. (Generalización)

Por el contrario, si partimos de elementos generales o premisas generales, del tipo «Todos los gatos tienen pelo» y tenemos un gato, teniendo en cuenta este tipo de razonamiento, por tanto, el gato tiene que tener pelo» entonces estamos hablando de razonamiento deductivo y como puede verse su rasgo principal es que parte de premisas generales para llegar a conclusiones particulares y por tanto a generar conocimiento.

Podemos, entonces, definir estos dos tipos de razonamiento, aunque grosso modo, de la siguiente manera:

- ▶ **Razonamiento inductivo:** parte de premisas particulares y obtiene conclusiones generales.
- ▶ **Razonamiento deductivo:** parte de premisas generales y obtiene conclusiones particulares.

Teniendo claro que el proceso de aprendizaje tiene entre sus componentes esenciales estos tipos de razonamiento, cabe preguntarse ¿Cómo entonces podría una máquina aprender?

Grosso modo, puede decirse que el campo del aprendizaje automático trata con los métodos que permiten dotar a los ordenadores de la capacidad de aprender. Establezcamos entonces los conceptos más importantes que nos permitan explicar cómo tiene lugar.

Supongamos que tenemos un programa de ordenador al que llamamos L y tenemos un profesor o entrenador al cual llamamos P y este, somete a L a un proceso de aprendizaje para ejecutar una tarea T , entregándole una experiencia D y además

Tema 1. Introducción al lenguaje automático

contamos con un método para evaluar lo que L ha aprendido y que llamaremos m que viene de métrica. Tanto L como P, se desenvuelven en un entorno llamado E.

Entonces, si L mejora su desempeño en la ejecución de la tarea T, dentro del entorno o ambiente que llamamos E, a medida que se le somete a la experiencia D, evaluado por P por medio de m, puede decirse que L, aprende.

Este proceso se muestra en la siguiente Figura.

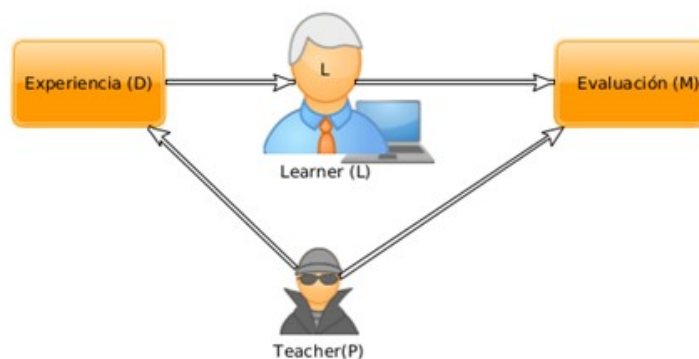


Figura 3. El proceso de aprendizaje grosso modo.

Pongamos como ejemplo, para entender mejor el proceso de aprendizaje automático, un programa de aprendizaje automático que debe aprender a jugar ajedrez.

En este caso concreto el aprendiz L debe aprender a jugar ajedrez, la cual es la tarea, a partir de una experiencia que tomaría la forma de partidas de ajedrez jugadas por grandes maestros. La métrica con la que evaluaríamos la mejora de L en la tarea T sería por medio de las jugadas ganadas, puede que jugando contra si mismos u otros L.

En los epígrafes siguientes iremos viendo la forma exacta que toman cada uno de estos elementos que hemos introducido de manera informal en este parte.

Tema 1. Introducción al lenguaje automático

Con estos elementos podemos, entonces, ver de manera también muy informal, los tipos de aprendizaje que pueden tener lugar.

Tipos de aprendizaje. Una aproximación informal

En un primer escenario, quizás el más elemental, pues es de los que utilizamos casi de manera intuitiva cuando educamos a nuestros hijos o por el que se entrenan las mascotas, es el siguiente:

Supongamos de nuevo que tenemos el aprendiz, llamado L y que interactúa de manera libre con el ambiente E, bajo la supervisión del profesor P y este lo remunera por las buenas acciones o lo castiga por las malas.

Este tipo de aprendizaje se le denomina aprendizaje con reforzamiento. Aunque este tipo de aprendizaje puede realizarse sin la supervisión de P y es E el que determina la remuneración o penalización a L. Este tipo de aprendizaje con reforzamiento se le llama por prueba y error.

Cuando el aprendizaje con reforzamiento se usa la remuneración de las buenas acciones entonces estamos haciendo un aprendizaje con reforzamiento positivo y viceversa cuando penalizamos las malas acciones hacemos uso del aprendizaje con reforzamiento negativo.

Esta es la manera en que amaestramos mascotas. Por ejemplo, si luego de indicar la ejecución de una acción, la mascota la ejecuta de manera correcta, pues se premia con un cubo de azúcar u otra acción. A esto se le llama reforzamiento positivo. Si por el contrario luego de ejecutar una acción fuera de las indicadas, como, por ejemplo, un perro que ladra a un transeúnte se tira de la correa y se regaña, se le da una reprimenda, se le penaliza esa conducta, entonces estamos en presencia de una forma de reforzamiento negativo.

Tema 1. Introducción al lenguaje automático

En este mismo escenario, donde tenemos, al mismo aprendiz que ya habíamos llamado L y un profesor llamado P y que actuaban en un ambiente E, también puede darse otro tipo de aprendizaje. En este caso, a L, el profesor P le entrega una lista de preguntas y sus respuestas, luego evalúa lo que L ha aprendido por medio de preguntas y puede con preguntas que no estaban en el cuestionario inicial. A este tipo de aprendizaje es al que se le llama aprendizaje supervisado.

- ▶ Aprendizaje supervisado

El detalle distintivo de este tipo de aprendizaje es precisamente que a L se le entregan las respuestas junto al cuestionario. Este L entonces aprende por medio de una experiencia ya obtenida y trata de responder preguntas que no están en el cuestionario.

También puede darse otra forma de aprendizaje, en la cual, a L, se le entregan solo las preguntas y luego P evalúa la manera en que L aprendió, es decir cómo se las arregló con las preguntas que se le hicieron. A este tipo de aprendizaje se le llama Aprendizaje no supervisado.

Y el rasgo distintivo en este tipo de aprendizaje es que, a L, no se les subministran las respuestas a las preguntas, él tiene que encontrarlas y eso es lo que se evalúa.

Para ir resumiendo, como hemos visto, en aprendizaje computacional podemos encontrarnos con tres tipos fundamentales de aprendizaje y digo fundamentales porque hay en la literatura algunos otros que no trataremos aquí.

- ▶ 1. Aprendizaje supervisado.
- ▶ 2. Aprendizaje no supervisado.
- ▶ 3. Aprendizaje con reforzamiento.

Tema 1. Introducción al lenguaje automático

Hemos establecido, de manera informal, cada uno de los tipos de aprendizaje y de la misma manera vamos a descubrir los elementos que conforman estos métodos y los acercaremos poco a poco al campo de ciencia de la computación.

Si analizamos detenidamente cada uno de estos métodos de aprendizaje, podemos darnos cuenta de que en todos ellos hay un elemento que se llama L , este representa el algoritmo o método de aprendizaje. Luego iremos viendo poco a poco que aspecto tiene y las formas que puede adoptar. También tenemos a un elemento P y este representa a los algoritmos de entrenamiento o los ingenieros de aprendizaje con sus herramientas que también poco a poco iremos poniéndole cara y cuerpo. Por otro lado, está, otro de los elementos importantes, la manera en que se le presenta a L la lista de cosas a aprender y que en la vida real toma forma de un conjunto de datos. Por último y no menos importante la manera en que T evalúa lo que L ha aprendido. Esta manera de evaluar lo que cada algoritmo aprende se denominan métricas y suelen ser variadas.

Por tanto, tenemos aquí, que en el aprendizaje computacional intervienen los siguientes elementos:

- ▶ Un conjunto de datos (la lista de preguntas con las respuestas o no).
- ▶ Un Algoritmo de aprendizaje (L en los ejemplos que pusimos).
- ▶ Una métrica de evaluación (en los ejemplos anteriores usamos m para denotarla)

1.2.3. El conjunto de datos, los tipos de variables y su clasificación

En esta parte estaremos aportando los detalles necesarios para comprender la forma que toma la experiencia, de la que hablamos en la sección anterior, que es la que se usa en el entrenamiento de los algoritmos de aprendizaje.

Esta experiencia de la que hemos hablado en las secciones anteriores toma la forma de Conjunto de Datos o Dataset. Específicamente a esta forma de experiencia se le

Tema 1. Introducción al lenguaje automático

llama Tabla de Decisión.

El conjunto de datos es el nombre que se le da, por tanto, a los datos que van a ser usados en el proceso de entrenamiento de los algoritmos de aprendizaje y están organizados de forma tabular, es decir, en filas y columnas al estilo de un libro de Excel.

Height	Weight	Age	Male
151.765	47.825606	63.0	1
139.700	36.485807	63.0	0
136.525	31.864838	65.0	0

Tabla 2. Tabla en Excel. Fuente: elaboración propia.

La Tabla anterior es un ejemplo de un *dataset* real, este contiene el alto, el peso, la edad y el sexo de 544 personas del pueblo Kung San, recopiladas por Nancy Howell en 1960. Este pueblo se encuentra ubicado al oeste del desierto de Kalahari entre Ovamboland y Botswana, al sur de África. Aquí solo se muestran tres casos; pero como ya se comentó este tiene 544 casos.

Un *dataset* está formado por filas y columnas. A las filas se les llama casos y aquí cada una de ellas representa a una persona a la que se les midió ciertas características (El peso, el alto, la edad y el sexo). Cada columna es una medición hecha a cada caso y se llaman variables, rasgos, predictores o características. Estos nombres se usan de forma equivalente en la literatura.

Estas variables o rasgos pueden clasificarse por el tipo de valor que asume o por la función que tienen a la hora de entrenar un modelo de aprendizaje, veamos en más detalles.

Si consideramos el uso de estas variables o rasgos, entonces pueden clasificarse en:

- Variables dependientes (variable predicha).

Tema 1. Introducción al lenguaje automático

- ▶ Variable Independientes (predictores).

Explicemos un poco esto para que se vea con claridad, si fuésemos a aprender a predecir el peso de una persona o caso del conjunto anterior y nos basamos en las otras variables que ya tenemos en él, como el alto, la edad y el sexo, entonces la variable dependiente es la queremos predecir, en este caso el peso y el resto de las variables son llamadas predictoras o variables independientes.

Si, por el contrario, queremos clasificar en hombre o no a los casos de este *dataset* entonces la variable *male*, sería la variable dependiente o predicha y el resto las variables independientes.

Ahora bien, estas variables también se pueden clasificar por su naturaleza, es decir por el tipo de valores que puede asumir y existen dos grandes grupos:

- ▶ Las cualitativas.
- ▶ Las cuantitativas.

A su vez, cada una de las anteriores se dividen dos:

Las variables cualitativas pueden dividirse en:

- ▶ Nominales.
- ▶ Ordinales.

Por otro lado, las cuantitativas se pueden dividir en:

- ▶ Discretas.
- ▶ Continuas.

Tema 1. Introducción al lenguaje automático



Figura 4. Clasificación de las variables según su naturaleza.

Ya sabemos que los cuatro tipos de variables son las siguientes:

- ▶ Nominales (categóricas).
- ▶ Ordinales.
- ▶ Discretas.
- ▶ Continuas.

Definamos y pongamos ejemplo de cada una de ellas.

Variables nominales (categóricas): estas se consideran como el nivel básico de medición. Esto entraña un proceso de codificación o un proceso por medio del cual se asignan números de forma arbitraria a las diferentes categorías que conforman la variable.

Las diferentes categorías simplemente constituyen una clasificación, ellas no pueden ordenarse de ninguna forma, son simplemente diferentes.

Tema 1. Introducción al lenguaje automático

Por ejemplo, cuando queremos preguntar a algún grupo de personas por qué razón trabajan, las razones pueden ser:

- ▶ Porque pagan bien.
- ▶ Porque el trabajo da una sensación de bienestar.
- ▶ Porque no hay mucha supervisión y se puede tomar sus propias decisiones.
- ▶ Porque el trabajo es excelente.

Por tanto, asignando un valor numérico a cada tipo de respuesta a esta pregunta estamos asignando una etiqueta a cada una de las posibles respuestas.

Hay veces que solo tenemos dos categorías Hembra/Varón, femenino/masculino, a estas además se les puede conocer como variables dicotómicas o binarias.

Este es el caso de la variable *male*, del *dataset* que analizamos como ejemplo, el de los pesos y el alto de 544 personas. En este caso la variable toma solo dos valores 0 y 1, para representar si el caso pertenece a un hombre o mujer.

Una de las características de las variables nominales es que carecen de orden.

Por tanto, una variable puede ser tratada como nominal cuando sus valores representan categorías que no obedecen a una clasificación intrínseca. Por ejemplo, el departamento de la compañía en el que trabaja un empleado. Algunos ejemplos de variables nominales son: región, código postal o filiación religiosa.

Otros ejemplos de variable nominal:

Outcome, en el *dataset* *pima-indians-diabetes-database*. Esta variable toma dos valores: 1 si el caso pertenece a un paciente con diabetes confirmada y 0 de lo contrario.

Tema 1. Introducción al lenguaje automático

Variables ordinales: este tipo de variable contiene valores que pueden ser ordenados. Por ejemplo, si le preguntamos a un grupo de trabajadores que escojan de entre las siguientes respuestas su nivel de satisfacción, entonces la variable puede considerarse ordinal:

Indique cuál de las siguientes respuestas representa su nivel de satisfacción en su trabajo.

- ▶ 1. Insatisfecho
- ▶ 2. Ni satisfecho ni insatisfecho
- ▶ 3. Medianamente satisfecho
- ▶ 4. Satisfecho
- ▶ 5. Muy satisfecho

Es evidente que aquí el número ofrece una medida de la cantidad de satisfacción con 1 siendo insatisfecho y 5 muy satisfecho. Eso es una escala.

Este tipo de variables es muy sencilla de identificar y no creo que sea necesario que pongamos muchos más ejemplos.

Variables discretas: son aquellas que toman valores puntuales en la recta numérica. Es decir, toman valores discretos.

Puede llamarse discreta a la variable Profesores, que contiene los números de profesores de Matemática en cada una de las universidades de España.

Otro ejemplo sería, la cantidad de niños que asisten a una clase y es importante hacer notar que, a la clase no pueden participar 8 niños y medio, por tanto, esta variable solo toma valores enteros.

Tema 1. Introducción al lenguaje automático

Otros ejemplos serían la cantidad de escalones de una escalera, el número de palabras en un texto, etc.

Variables Continuas: Este tipo de variables pueden tomar cualquier valor numérico. Es decir, estos valores pueden ser enteros o no.

Por ejemplo, la altura de todos los estudiantes de la escuela x, o la longitud de las barras de acero salidas de una planta, o el ingreso de los trabajadores de una fábrica.

Otro ejemplo de este tipo de variable sería la distancia recorrida por un vehículo al frenarlo a una determinada velocidad. O el Valor del Tiempo de Vida de un consumidor de una gran cadena de supermercados.

En el caso del *dataset* que pusimos de ejemplo, las variables *height* y *weight*, son variables continuas.

Los algoritmos y los modelos de aprendizaje automático

Ya sabemos, de los epígrafes anteriores, que se entiende por Inteligencia Artificial débil a aquellos métodos que se concentran en solucionar problemas puntuales de la vida real, como poder clasificar las variantes de covid-19 a partir de su perfil genético o a clasificar imágenes de personas en sistemas de alarmas que usan IA.

También, sabemos que el aprendizaje automático se desarrolló como una herramienta para resolver los problemas anteriores y por tanto es uno de los métodos de la IA débil.

Sabemos también que se puede aprender por medio del razonamiento inductivo o deductivo y que el primero se basa en premisas particulares, específicas y puede llegar a conclusiones generales. Es muy importante que tomemos nota de este detalle pues todos los algoritmos de aprendizaje automático tienen este mecanismo de razonamiento y al cual se le llama principio inductivo.

Tema 1. Introducción al lenguaje automático

Es, este principio inductivo la forma práctica o computacional que toma el razonamiento o inferencia inductivos cuando es aplicado al aprendizaje automático.

En esta sección veremos, aunque también de manera muy informal, cuáles son los algoritmos más usados para cada uno de los métodos de aprendizaje que vimos en las secciones anteriores y estableceremos la diferencia entre algoritmo y modelo.

También sabemos, de manera más o menos clara la forma que toma la experiencia a partir de la cual los algoritmos pueden aprender; pero es necesario hacer algunas distinciones en lo que respecta a las tareas particulares que se pueden llevar a cabo dentro de cada uno de los tipos de aprendizaje.

¿A qué me refiero con esto?, pues, que en dependencia del tipo tarea que se pretende ejecutar así es el tipo de algoritmo que debemos escoger o tenemos disponible para usar.

Si queremos, por ejemplo, predecir al precio de una casa a partir del lugar donde está ubicada y los metros cuadrados, etc. Debemos tener en cuenta que la variable precio, debe ser una del tipo continua y eso quiere decir que puede tomar cualquier valor numérico en el conjunto de los reales positivos.

Por tanto, cuando tenemos esta situación, donde la tarea requiere predecir una variable continua, estamos en presencia de una tarea de Regresión y por tanto hay que escoger algoritmos de regresión.

Si, por el contrario, se necesita predecir el género de una persona a partir de rasgos como el peso y el alto, o se necesita predecir si un paciente padece cardiopatía hipertensiva a partir de determinados factores de riesgo, entonces la variable género, o la variable HTA son variables categóricas dicotómicas, entonces estamos ante una tarea de Clasificación.

Tema 1. Introducción al lenguaje automático

Entonces, podemos decir que en el aprendizaje supervisado se pueden acometer dos tipos de tareas básicamente:

Aprendizaje supervisado:

- ▶ Regresión.
- ▶ Clasificación.

Es necesario tener en cuenta, aunque por ahora a modo de información, que hay algoritmos de clasificación que en vez de retornar una variable dicotómica tal como Hombre/mujer, 0/1, retornan un valor que corresponde a la probabilidad de que un caso n pertenezca a una clase c . Es decir, el valor que obtendríamos sería 0,78, por ejemplo, o 0,51.

Si a diferencia de lo que pasa en aprendizaje supervisado, donde contamos con una tabla de datos que contiene una variable de decisión que ayuda a los algoritmos a aprender la dependencia de esta a partir de los predictores, tenemos una tabla con n casos con b variables y tenemos que aprender de ellos, entonces estamos ante una tarea de *clustering* o agrupamiento. Este tipo de tarea pertenece a los métodos de aprendizaje no supervisado.

En el agrupamiento (*clustering*) se utiliza una medida de distancia o similitud entre los elementos del conjunto para agruparlos. La idea es formar *clusters* con los elementos cuya similitud sea mayor que los de otros grupos.

Por tanto, el elemento a destacar de este tipo de aprendizaje son las medidas de similitud o de distancia que usamos. El método usado por excelencia en aprendizaje no supervisado es el de agrupamiento o *clustering* y este también puede dividirse de dos formas: jerárquico y particional.

Tema 1. Introducción al lenguaje automático

Aprendizaje no supervisado:

► *Clustering:*

- Jerárquico.
- Particional.

Ahora bien, el aprendizaje con reforzamiento es un paradigma completamente diferente a los casos de aprendizaje supervisado y no supervisado. A diferencia del caso supervisado que funciona bajo la idea de la inferencia inductiva, el aprendizaje con reforzamiento se basa en la idea del llamado «condicionamiento operante». Este principio, hablando grosso modo, sostiene la idea de que, si después de cada acción se ejecuta otra en forma de recompensa o castigo, se puede modificar la frecuencia de ejecución de las acciones anteriores por medio del refuerzo o debilitamiento de la motivación a ejecutarla y por tanto nos lleva a aprender.

A diferencia del aprendizaje supervisado, en este, solo existe una señal de recompensa. Es también, importante saber que en este tipo de aprendizaje la recompensa puede ser inmediata o demorada.

Para atrapar la idea de aprendizaje con reforzamiento, es necesario verlo como un proceso secuencial, donde en cada paso t , se observa el estado y del ambiente y se elige una acción a y se ejecuta y como resultado se obtiene una recompensa R . Este proceso se ejecuta hasta que se llega a un objetivo/meta.

Tema 1. Introducción al lenguaje automático

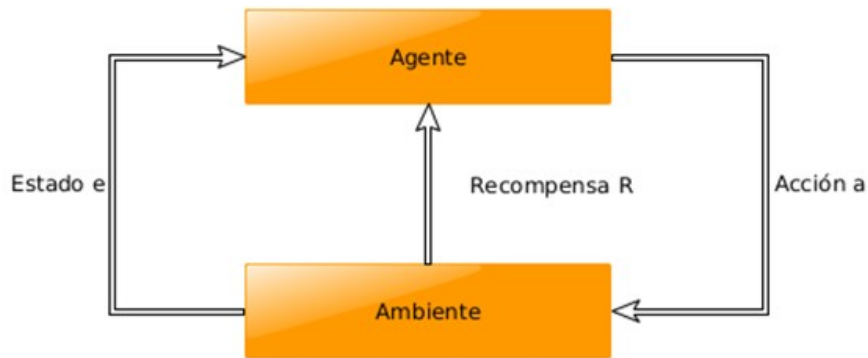


Figura 5. Aprendizaje con reforzamiento.

Aunque hay otros algoritmos de aprendizaje con reforzamiento, los dos más usados son Q-Learning y SARSA.

Entonces, a modo de conclusión de esta parte haremos un listado de los algoritmos disponibles para cada uno de los tipos de aprendizaje. Es necesario explicar que esta lista no es completa, ni pretende serlo, solo daremos los más usados o los que con más frecuencia se pueden encontrar en la literatura.

Algoritmos usados en aprendizaje supervisado:

► Algoritmos de regresión:

- Regresión lineal
- Regresión logística.

► Algoritmos de clasificación:

- KNN. (Algoritmo de los K vecinos más cercanos: k-Nearest Neighbor).
- CART (Arboles de clasificación y regresión).
- Random Forest. (Lo dejo en inglés a falta de un buen término para la su traducción).

Tema 1. Introducción al lenguaje automático

- Naives bayes.
- Máquinas con vectores de apoyo. (Support Vector Machines).
- Redes Neuronales. (Hay muchos modelos englobados en este término).

Los algoritmos de agrupamiento pueden dividirse a grandes rasgos en tres:

- ▶ Particionales.
- ▶ Jerárquicos.
- ▶ Basados en densidad.

Sin embargo, por su simpleza y uso solo tocaremos uno. Pero listaremos algunos.

Algoritmo de Clustering:

- ▶ K-Mean.
- ▶ C-Mean.
- ▶ Fuzzy C-Mean.
- ▶ DBSCAN.

Algoritmos de aprendizaje con reforzamiento:

- ▶ Q-Learning.
- ▶ SARSA.

Creo que sería muy útil poner algunos ejemplos de problemas en los que se puede aplicar las técnicas de aprendizaje automático y ver cómo se seleccionan basados en el tipo de problema.

Ejemplo 1

Tema 1. Introducción al lenguaje automático

Supongamos que somos parte de un equipo de análisis de datos de la empresa x y la dirección de Marketing nos encomienda la tarea de desarrollar un modelo para la predicción del CLV (Valor del tiempo de vida de un consumidor) a partir de los datos históricos almacenados.

Pues si ya tenemos claro que la tarea es la de predecir el Valor del tiempo de vida de un cliente, pero no sabemos cómo luce el modelo analítico que se usa para calcularlo, pero sabemos que contamos con varias mediciones que intervienen de manera directa en su valor, pues procedemos a juntar dichos rasgos y evaluar los datos con que contamos.

Inicialmente ya sabemos que CLV es una variable continua y la unidad de medida es el Euro. Esto nos dice claramente que como vamos a predecir una variable continua tenemos que usar algoritmos de regresión.

Entonces, lo que hay que sacar de este ejemplo es que si se predicen variables continuas los algoritmos que se pueden utilizar son los de regresión.

Ejemplo 2

Como en el ejemplo anterior, que ya formábamos parte del equipo de análisis de la empresa X, la dirección de marketing ahora nos propone otro problema, automatizar la tarea de segmentación de los clientes usando la técnica llamada RFM.

Bueno, al iniciar el análisis del problema, nos referimos como siempre a los datos con los que contamos para la solución y nos damos cuenta que los datos que tenemos para eso son la cantidad en dinero que ha gastado un cliente, la frecuencia con la que el cliente compra nuestro producto y la Recencia o el tiempo que media entre una compra y otra. En este tipo de

Tema 1. Introducción al lenguaje automático

problemas, no tenemos una variable o rasgo que nos permita aprender una dependencia funcional y por tanto tenemos que encontrar un método que aprenda de la propia estructura de los datos.

Cuando tenemos este tipo de problemas, estamos en presencia de un problema de Clustering o agrupamiento y podemos atacarlo usando k-mean u otro de los algoritmos.

Ejemplo 3

Supongamos, para este ejemplo, que somos parte de una empresa consultora Y y se nos acerca un cliente del sistema de salud interesado en un sistema que mejore las propiedades de los sistemas de ventilación mecánica para la atención a pacientes críticos de Covid-19. La mayoría de estos sistemas son controlados por controladores PID, pero sería posible desarrollar un sistema controlado por medio de Aprendizaje con reforzamiento y que aprenda por medio de miles de simulaciones. En este tipo de sistema de control el Learner aprendería por medio de la señal de recompensa cuando lo está haciendo bien o mal a partir de los signos asociados al sistema respiratorio.

Hasta aquí hemos visto los nombres de algunos de los algoritmos que se usan en aprendizaje automático. Hemos visto también algunos ejemplos de los tipos de problemas que pueden resolverse usando regresión/clasificación, clustering y aprendizaje con reforzamiento. Sin embargo, hay un elemento que no hemos mencionado y que es sumamente importante que se entienda desde el principio. En la literatura de aprendizaje automático se trata con muchísima frecuencia, de manera intercambiable los términos Modelo y algoritmo, lo que resulta confuso para los que se inician.

Tema 1. Introducción al lenguaje automático

Es importante entonces que dejemos claro que es un algoritmo de aprendizaje automático y que es un modelo.

Grosso modo, un algoritmo es una serie de pasos que se utilizan para la determinación o cálculo de los parámetros de un modelo a partir del conjunto de datos. Mientras el algoritmo es el mismo para todos los conjuntos de datos, el modelo puede ser distinto. Por tanto, de manera muy general podemos afirmar que el algoritmo es el método por el cual adquirimos el conocimiento para ejecutar una tarea, mientras que el modelo es el conocimiento en sí, el que nos permite luego responder preguntas.

Para ejemplificar la diferencia entre Algoritmo y modelo, tomemos el caso de una red neuronal. No importa que no entienda ahora todo lo relacionado con el algoritmo de Redes Neuronales.

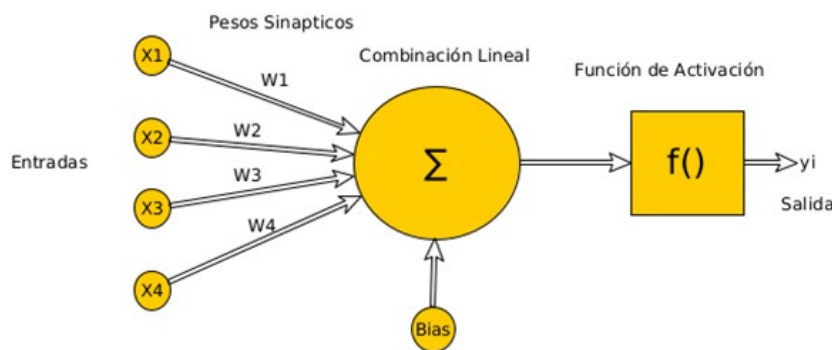


Figura 6. Esquema simplificado de una Red Neuronal.

Cuando usamos una red neuronal para la solución de algún problema de aprendizaje automático, al algoritmo se le pasa el conjunto de datos y este inicia un proceso que se denomina entrenamiento. Este «entrenamiento» no es otra cosa que una serie de pasos para calcular la matriz de pesos sinápticos w , a partir de cada una de los casos del conjunto de datos. Una vez calculada esta matriz de pesos podemos decir que tenemos un modelo entrenado. Es a esta matriz junto a otros metadatos a lo que

Tema 1. Introducción al lenguaje automático

se puede llamar modelo y ser guardado o desplegado a producción.

En la imagen, los círculos con nombre $x_1...x_4$ son los rasgos o variables del conjunto de datos y los $w_1...w_4$ son los llamados pesos sinápticos. A $f()$ se le llama función de activación. El resto de los detalles de cómo funciona este tipo de algoritmo se dará en la parte correspondiente a los detalles de cada uno de los modelos.

Podemos definir entonces:

Algoritmo de aprendizaje automático es el conjunto de pasos necesarios para el cálculo de los parámetros del modelo a partir del conjunto de datos de entrada.

Modelo de aprendizaje automático es el conjunto de parámetros arrojado por algún algoritmo de aprendizaje después del entrenamiento. Este nos permite hacer predicciones.

El proceso de aprendizaje automático

Utilizar un algoritmo de aprendizaje automático implica de alguna forma un conjunto de tareas que se ordenan de forma natural y que facilitan la obtención de modelos y su uso.

Hay autores que dividen este proceso hasta en 7 etapas, nosotros lo dejaremos en cuatro que incluyen todos los pasos necesarios para el uso del aprendizaje automático.

Por tanto, las etapas que componen el proceso de uso del aprendizaje automático son:

- ▶ Establecimiento del problema.
- ▶ Recolección y limpieza de los datos.
- ▶ Construcción y evaluación de los modelos.

Tema 1. Introducción al lenguaje automático

- Despliegue.

Etapas del proceso de aprendizaje automático:

Etapas de establecimiento del problema:

En esta etapa se comienza a tratar de entender el problema que necesitamos resolver. Esto implica, muchas veces, reuniones frecuentes con los que saben del problema. Si tomamos como ejemplo los casos que ya vimos en el epígrafe anterior, pues deberíamos tener sesiones interminables con la gente de Marketing para entender de que trata el problema de la predicción del valor del tiempo de vida del cliente y revisar los métodos por los que antes se hacía y leer documentación, buscar experiencias de otras empresas. En fin tratar de entender que se nos pide.

Etapas de recolección y limpieza de los datos:

Una vez que entendemos que es lo que nos piden, podemos pasar a evaluar los datos que tenemos y si no los tenemos pues tratar de encontrarlos. Luego que ya recopilamos los datos que pensamos necesitamos, iniciamos una etapa de análisis exploratorio, donde tratamos de entender que nos cuentan los datos. Revisar si faltan valores, si los datos tienen anomalías o están en el formato correcto para que los algoritmos puedan procesarlos.

Esta es una de las etapas a las que normalmente se le dedica más tiempo.

Es sumamente importante que se entienda desde el inicio que no todos los algoritmos pueden funcionar con todos los tipos de variables que puede contener los conjuntos de datos.

También en esta etapa hay que ejecutar varias acciones sobre los datos que de no hacerse provocarían que los modelos no pudieran operar de forma correcta. Entre ellas se encuentra el balanceo de los casos, la selección de las variables que no son necesarias, la normalización, el centrado, etc. Estas operaciones las trataremos en