

Título del trabajo (elegir uno original)

Domínguez-Adame Ruiz, Alberto

Dpto. Ciencias de la Computación e Inteligencia Artificial
Universidad de Sevilla
Sevilla, España
albdomrui@alum.us.es

Vilaplana de Trias, Francisco David

Dpto. Ciencias de la Computación e Inteligencia Artificial
Universidad de Sevilla
Sevilla, España
fravilde1@alum.us.es

Resumen—Escribir aquí dos párrafos indicando el objetivo principal del trabajo, y un resumen de las conclusiones obtenidas. Cabe mencionar que este documento se ha confeccionado siguiendo el formato de conferencias de IEEE (ver la guía para autores para más información, existen plantillas para Word y L^AT_EX). Este documento se debe emplear como guía, se pueden añadir nuevas secciones según sea necesario. Es importante dotarlo de un número razonable de referencias bibliográficas.

Palabras clave—Inteligencia Artificial, otras palabras clave...

I. INTRODUCCIÓN

Para empezar, el aprendizaje se define como la adquisición del conocimiento de algo por medio del estudio, el ejercicio o la experiencia, en especial de los conocimientos necesarios para aprender algún arte u oficio. Al hablar sobre el aprendizaje automático estaríamos entrando en el campo de la Inteligencia Artificial (IA). Podemos definir esta como un programa de computación diseñado para realizar determinadas operaciones que se consideran propias de la inteligencia humana, como el Aprendizaje Automático (Machine Learning), una rama de la inteligencia artificial, cuyo objetivo es el desarrollo de un sistema (modelo matemático que realiza una determinada tarea) el cual tiene un mejor desempeño con la experiencia, dados una serie de datos.

Nuestro estudio hace uso de datos relacionales, los cuales están unidos entre sí (tienen una relación) que queda representado como aristas en un grafo, entendiendo un grafo como un conjunto de nodos unido (o no) mediante aristas. En la Fig. 1 podemos ver un ejemplo de un grafo a partir de unos datos relacionales.

Para la realización de este trabajo se ha hecho uso del entorno de desarrollo Jupyter y de la herramienta notebook. El código está escrito mediante el lenguaje de programación Python, usando principalmente las librerías pandas, NetworkX y scikit-learn, y otras como keras, tensorflow y matplotlib.

Los modelos de clasificación usados son: KNN, Naive Bayes, Árboles de Decisión y Redes Neuronales.

II. CONJUNTO DE DATOS Y TAREA DE PREDICCIÓN A REALIZAR

Hemos elegido un dataset (fíjese en la Fig. 2) de la página de streamings en directo twitch.tv, en la que se ve: las visitas, días (ambos de tipo entero) que transmitido un canal en directo, si tiene contrato o no con twitch (partner de tipo Boolean) y si

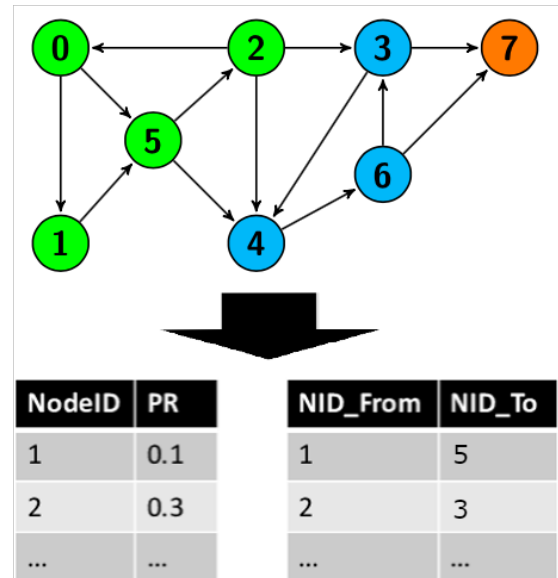


Fig. 1. Ejemplo de Grafo

el contenido del canal (id de tipo Integer) es o no para adultos (mature de tipo Boolean). Al ser la propia Twitch la que ofrece este contrato, hemos decidido que la predicción para nuestros modelos sea el atributo partner (true/false). Las aristas (edges) representan si dos usuarios tienen una relación de amistad.

Este dataset tenía dos tipos de id asociados a cada usuario, uno generado aleatoriamente y otro donde al usuario se le asignaba un valor entre 0 y el numero total de usuarios que recoge el conjunto de datos, por lo que decidimos que eliminar el id generado de forma aleatoria para mayor claridad y sólo hace falta un id por usuario para identificarlos.

Para este estudio hemos decidido que el atributo a predecir será partner. Para ello, haremos uso de los atributos: days, views y mature (el id no es relevante en este caso). Lo hemos planteado como una tarea de clasificación binaria, ya que partner solo tiene dos valores posibles (True o False).

Por último hemos codificado (codificación one-hot) los valores de mature de Booleano a Integer, pasando True-False a 1-0 respectivamente, para que sea más eficiente estimar las probabilidades del atributo partner.

	id	partner	days	views	mature
0	2299	False	1459	9528	0
1	153	False	1629	3615	1
2	397	False	411	46546	1
3	5623	False	953	5863	1
4	5875	False	741	5594	1
...
7121	3794	False	2624	3174	0
7122	6534	False	2035	3158	1
7123	2041	False	1418	3839	1
7124	6870	False	2046	6208	1
7125	3919	False	1797	3545	0

Fig. 2. Tabla de Datos

III. MÉTRICAS RELACIONALES

Como este es un trabajo de desarrollo de modelos orientados al aprendizaje automático grafos, hemos elegido una serie de métricas relacionales que aplicaremos como atributos relacionales. En concreto hemos elegido tres métricas:

- 1) **Betweenness centrality** (Centralidad de intermediación): La centralidad en un grafo puede ser entendida como una medida que determina la relevancia de un nodo dentro del grafo y permite comparar o contrastar dicho vértice con otros. La Betweenness Centrality es una medida de centralidad que cuantifica la frecuencia en la que un nodo se encuentra en el camino más corto entre dos nodos determinados. Cuando en un grafo existen nodos de alta intermediación, estos suelen jugar un rol importante en la estructura a la que pertenecen. Estos nodos también poseen capacidades de ser controladores o reguladores de los flujos de información dentro de la estructura total del grafo.
- 2) **Clustering** (agrupamiento): es una tarea que tiene como finalidad principal lograr el agrupamiento de nodos, para lograr construir subconjuntos de datos conocidos como Clusters.
- 3) **Degree** (grado): Número de aristas (relaciones) que inciden sobre el nodo.

IV. ALGORITMOS DE APRENDIZAJE AUTOMÁTICO

A continuación exponemos los algoritmos que hemos utilizado con nuestro dataset:

A. KNN

El algoritmo kNN (del inglés *k Nearest Neighbors*, *k* vecinos más cercanos), es un clasificador de aprendizaje supervisado no paramétrico, por lo que la cantidad de parámetros del modelo coincide exactamente con la cantidad de ejemplos de entrenamiento. Se basa en la proximidad de los *k* nodos más cercanos calculando la cercanía a partir de los atributos, si estos fueran discretos sería necesaria una codificación de los mismo.

Analizando nuestro dataset, se puede observar como el atributo "views" toma valores de mayor magnitud con respecto al resto lo que provoca que tenga más en cuenta para la predicción del modelo. Con el fin de evitar esta situación se hace uso de la normalización a los atributos del conjunto de entrenamiento. Existe una gran variedad de maneras de normalizar un atributo numérico, nosotros hemos optado por el método min-max [1] que consiste en que a cada atributo restarle el valor mínimo de este y dividir esta diferencia entre la resta del valor máximo y mínimo:

$$v'_i = \frac{v_i - m}{M - m} \quad (1)$$

Para realizar kNN es necesario dar un valor al hiperparámetro *k*, referido al número de vecinos que se va a tener en cuenta a la hora de clasificar un nuevo dato, nosotros le hemos dado valores del uno al diez para más tarde evaluar el rendimiento de cada uno y elegir el mejor. También existen métodos para calcular la distancia, en nuestro caso hemos usado:

- Distancia de *Hamming* para el dataset sin atributos relacionales:

$$d(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^n \mathbb{1}(x_i \neq x'_i) \quad (2)$$

Hemos usado esta métrica ya que según la documentación de la librería [2] es la recomendada para valores enteros.

- Distancia *Manhattan* para el dataset con atributos relacionales:

$$d(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^n |x_i - x'_i| \quad (3)$$

En este caso es la adecuada al tratarse de valores numéricos reales.

Para el cálculo del modelo ha sido necesaria la clase *neighbors* de la librería *sklearn*. Por último, una vez se ha construido el modelo lo evaluamos mediante validación cruzada (cross validation) con 10 pliegues. Se divide el dataset según el número de pliegues, siendo uno de estos el que se usará para probar el modelo y el resto para entrenarlo, de tal manera que cada pliegue acabe siendo usado como subconjunto de prueba. Cada una de estas iteraciones dará como resultado una tasa de acierto, la media de estas será el valor asociado a cada *k*. Por último se comparan y se escoge el valor de *k* que obtuviera la media la más alta.

B. Naive Bayes

C. Árboles de Decisión

Usa los comandos `table` y `tabular` para iniciar una tabla simple — mira la tabla I, como ejemplo.

V. RESULTADOS

En esta sección se detallarán tanto los experimentos realizados como los resultados conseguidos:

l para left	c para centro	r para derecha
Ejemplo	Centrado	Alineado a la
Izquierda	13	Derecha

TABLA I
UNA SIMPLE TABLA.

TABLA II
EJEMPLO DE TABLA

A	B	C
1	2	3
4	5	6

- Los experimentos realizados, indicando razonadamente la configuración empleada, qué se quiere determinar, y como se ha medido.
- Los resultados obtenidos en cada experimento, explicando en cada caso lo que se ha conseguido.
- Análisis de los resultados, haciendo comparativas y obteniendo conclusiones.

Se pueden hacer uso de tablas, como el ejemplo de la tabla II.

VI. CONCLUSIONES

Finalmente, se dedica la última sección para indicar las conclusiones obtenidas del trabajo. Se puede dedicar un párrafo para realizar un resumen sucinto del trabajo, con los experimentos y resultados. Seguidamente, uno o dos párrafos con conclusiones. Se suele dedicar un párrafo final con ideas de mejora y trabajo futuro.

VII. BIBLIOGRAFÍA

REFERENCIAS

- [1] <https://scikit-learn.org/stable/modules/preprocessing.html#normalization>
- [2] <https://scikit-learn.org/stable/modules/generated/sklearn.metrics.DistanceMetric.html>