

Práctica Netflix

Tecnologías para el Análisis de Datos Masivos

Daniel Ramos & Sergi Fornés

Importación y limpieza de los datos

Información sobre el Raw data

Vamos a realizar el análisis a partir de datos sobre valoraciones de películas de Netflix. Los datos que nos pueden resultar más interesantes son las notas de las valoraciones y la identificación de quien hace estas valoraciones.

Obtenemos los datos a partir de 5 ficheros `.txt` y un `.csv`.

`combined_data_x.txt`

Tenemos a nuestra disposición 4 ficheros de este tipo: `combined_data_1.txt`, `combined_data_2.txt`, `combined_data_3.txt` y `combined_data_4.txt`. En estos ficheros tenemos información sobre valoraciones numéricas puestas por usuarios a películas de Netflix. Cada bloque de valoraciones está precedido por un número que indica la película (ID de la película). Para cada película existe un conjunto de valoraciones, de las cuales tenemos información sobre la fecha de la valoración, la nota (del 1 al 5) y un identificador del usuario.

`filas_ID_combined_all.txt`

En el fichero tenemos la localización de los bloques de valoraciones por películas en los ficheros `combined_data_x.txt`.

Columnas del fichero:

- `X1`: Identificación del bloque. Es el mismo nombre que aparece en los ficheros `combined_data_x.txt`.
- `fila`: Número de fila en el que se encuentra la identificación del bloque.
- `ID`: ID de la película asociada al bloque.
- `fila_final`: Última fila que contiene una valoración de la película `ID`.
- `data`: Número de fichero `combined_data_x.txt` al que pertenece la película `ID`. Tiene 4 valores posibles.

`movie_titles.csv`

En este archivo podemos encontrar información sobre las películas. Tenemos el ID de la película, su nombre en inglés y su año de estreno.

Obtención de las películas de interés

Creación del data frame

Análisis exploratorio de los datos