

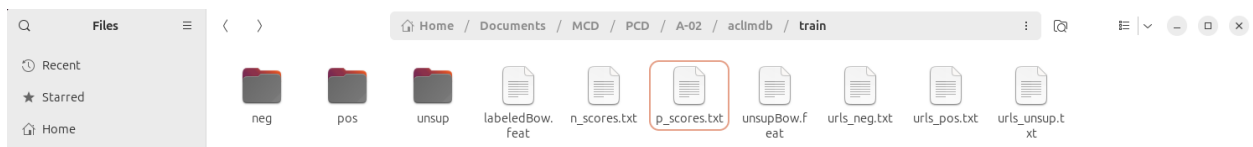
Programación para la Ciencia de Datos
Eddel Elí Ojeda Avilés
17 de agosto, 2024

Para encontrar las 10 películas peor y mejor calificadas en promedio del directorio train primero nos posicionamos en dicho directorio.

```
eddel@eddel-VirtualBox: ~/Documents/MCD/PCD/A-02/acLimdb/train
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$
```

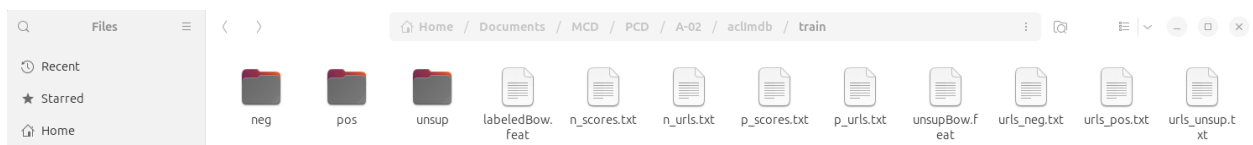
Posteriormente a partir de la lista archivos en los directorios neg y pos guardamos únicamente la puntuación asignada en las reseñas en los archivos n_scores.txt y p_scores.txt, respectivamente. Esto se hizo usando el comando awk para tomar el segundo elemento en los títulos de los archivos, una vez señalados los separadores “_” y “.”.

```
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ cd Documents/MCD/PCD/A-02/acLimdb/train
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ ls neg | awk -F'[_.]' '{print $2}' > n_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ ls pos | awk -F'[_.]' '{print $2}' > p_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$
```



Con el fin de utilizar las urls apropiadas, creamos los archivos n_urls.txt y p_urls.txt a partir de los archivos urls_neg.txt y urls_pos.txt, respectivamente, los cuales contienen las urls resultantes de eliminar “/usercomments” con el comando sed.

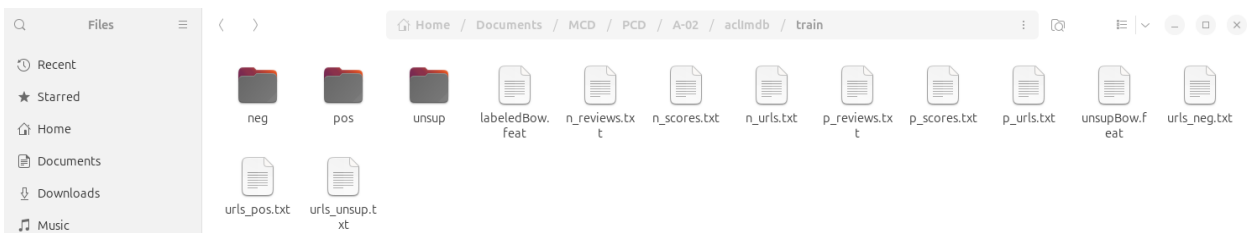
```
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ cd Documents/MCD/PCD/A-02/acLimdb/train
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ ls neg | awk -F'[_.]' '{print $2}' > n_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ ls pos | awk -F'[_.]' '{print $2}' > p_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ sed 's/\/usercomments$//' urls_neg.txt > n_urls.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ sed 's/\/usercomments$//' urls_pos.txt > p_urls.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$
```



La idea a continuación sería crear dos archivos n_names.txt y p_names.txt con los títulos de las películas asociadas a cada una de las urls, esto usando el comando curl, sin embargo, sigo trabajando en eso.

El siguiente paso consiste en utilizar el comando paste para generar los archivos n_reviews.txt y p_reviews.txt mezclando los archivos que contienen las urls, títulos de películas y puntuaciones negativas y positivas, respectivamente.

```
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ cd Documents/MCD/PCD/A-02/acLimdb/train
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ ls neg | awk -F'[_.]' '{print $2}' > n_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ ls pos | awk -F'[_.]' '{print $2}' > p_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ sed 's/\/usercomments$//' urls_neg.txt > n_urls.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ sed 's/\/usercomments$//' urls_pos.txt > p_urls.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ paste n_urls.txt n_scores.txt > n_reviews.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$ paste p_urls.txt p_scores.txt > p_reviews.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acLimdb/train$
```

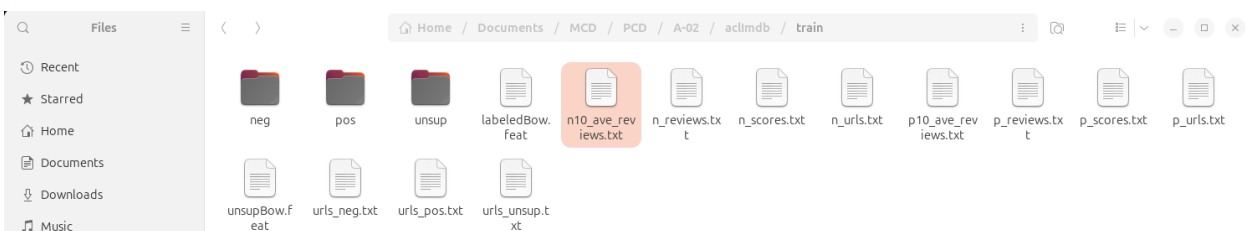


Finalmente, creamos los archivos `n10_ave_revies.txt` y `p10_ave_revies.txt` usando el comando `awk` donde calculamos la suma de las puntuaciones asociadas a la misma url a la par de que utilizamos un contador para guardar el número de veces que cada url se repite, para posteriormente guardar en un archivo cada una de las url distintas seguidas del cociente de la suma de sus puntuaciones entre cantidad de reseñas, es decir, su puntuación promedio, ordenándolas por su segunda componente de manera descendente o ascendente según el tipo de reseñas a considerar, para finalmente tomar únicamente las primeras o últimas 10 dependiendo de si se buscan los promedios más pequeños o más grandes.

```

eddel@eddel-VirtualBox: ~/Documents/MCD/PCD/A-02/acldmb/train
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ ls neg | awk -F'[_ .]' '{print $2}' > n_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ ls pos | awk -F'[_ .]' '{print $2}' > p_scores.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ sed 's/\\/usercomments$//' urls_neg.txt > n_urls.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ sed 's/\\/usercomments$//' urls_pos.txt > p_urls.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ paste n_urls.txt n_scores.txt > n_reviews.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ paste p_urls.txt p_scores.txt > p_reviews.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ awk '{n_reviews[$1]+=$2; count[$1]++} END {for (url in n_reviews) print url, n_reviews[url]/count[url]}' n_reviews.txt | sort -k2 -n | head -n 10 > n10_ave_reviews.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$ awk '{p_reviews[$1]+=$2; count[$1]++} END {for (url in p_reviews) print url, p_reviews[url]/count[url]}' p_reviews.txt | sort -k2 -n | tail -n 10 > p10_ave_reviews.txt
eddel@eddel-VirtualBox:~/Documents/MCD/PCD/A-02/acldmb/train$

```



```

n10_ave_reviews.txt
http://www.imdb.com/title/tt0000439 1
http://www.imdb.com/title/tt0001032 1
http://www.imdb.com/title/tt0007558 1
http://www.imdb.com/title/tt0008458 1
http://www.imdb.com/title/tt0009123 1
http://www.imdb.com/title/tt0011588 1
http://www.imdb.com/title/tt0013422 1
http://www.imdb.com/title/tt0016104 1
http://www.imdb.com/title/tt0016392 1
http://www.imdb.com/title/tt0018328 1

```

```

p10_ave_reviews.txt
http://www.imdb.com/title/tt0429203 10
http://www.imdb.com/title/tt0430650 10
http://www.imdb.com/title/tt0432028 10
http://www.imdb.com/title/tt0443678 10
http://www.imdb.com/title/tt0495047 10
http://www.imdb.com/title/tt0688816 10
http://www.imdb.com/title/tt0751121 10
http://www.imdb.com/title/tt0770745 10
http://www.imdb.com/title/tt0815744 10
http://www.imdb.com/title/tt0961088 10

```