



Informe Movie Data

Proyecto Analítico de Datos

Profesor

Mario Ortiz

Alumno

Johan Órdenes

Leonel Villagra

Fecha

27/10/2020

Índice

Análisis de los datos	3
Limpieza de los datos	4
Comparación de datos	6
Tiempo de compilación	6
Exportación de archivos	7
Carga a la bases de datos	8

En esta investigación se nos pide hacer el proceso de ETL, es decir exportar, transformar y cargar los datos del archivo moviedata.xlsx. Esto lo hicimos realizando la lectura de los datos entregados, limpieza y eliminación de los datos inservibles, y por último el almacenamiento de los datos en una base de datos en Postgres. Para esto separamos nuestro código en dos fases. La primera para la lectura y guardado de archivos que veremos más adelante, y la segunda para subir netamente dichos archivos a la base de datos.

1. Análisis de los datos

color	director_name	num_critics_for_reviews	duration	director_facebook_likes	actor_3_facebook_likes	actor_2_name	actor_1_facebook_likes	gross
Color	Mikael Håfström	286	115	101	585	50 Cent	13000	25121291
Color	Brian A Miller	46	93	32	1000	50 Cent	13000	
Color	Antoine Fuqua	305	124	845	424	50 Cent	15000	52418902
Color	Don Coscarelli	81	97	152	135	A. Michael Baldwin	674	7282851
Color	Randall Rubin	23	99	0	131	A.J. Buckley	363	
Color	Bobby Roth	1	120	40	249	Aaron Ashmore	912	
Color	Alejandro Amenábar	140	106	448	648	Aaron Ashmore	9000	54734
Color	John Dahl	121	93	131	90	Aaron Hughes	497	2426851
Color	Pou-Soi Cheang	14	119	3	22	Aaron Kwok	879	
Color	Gary Winick	91	78	56	184	Aaron Stanford	376	2882062
Color	Kirsten Sheridan	148	114	59	403	Aaron Staton	49000	31655091
Color	D.J. Caruso	253	105	154	501	Aaron Yoo	884	80050171
Color	D.J. Caruso	253	105	154	502	Aaron Yoo	884	80050171
Color	Jonathan Levine	147	99	129	362	Aaron Yoo	976	2077046
Color	Larry Charles	313	99	119	174	Asif Mandvi	600	59617068
Color	Shekhar Kapur	202	114	159	877	Abbie Cornish	13000	16264475
Color	Zack Snyder	435	128	0	826	Abbie Cornish	4000	36381716
Color	José Padilha	492	117	294	1000	Abbie Cornish	10000	58607007
Color	José Padilha	492	117	294	1000	Abbie Cornish	10000	58607007
Color	R. Balki	12	133	12	106	Abhishek Bachchan	464	199228
Color	Michael Apted	116	116	150	956	Abigail Spencer	18000	6002756
Color	McG	241	103	368	520	Abigail Spencer	27000	54758461
Color	Shawn Levy	156	103	189	949	Abigail Spencer	2000	34290142
Color	Sam Weisman	53	92	39	458	Abraham Benrubi	3000	105263257
Color	Craig Gillespie	178	117	44	531	Abraham Benrubi	788	27550735
Color	Eric England	15	81	15	93	Ace Marrero	847	
Color	Frank Oz	140	85	0	490	Adam Alexi-Malle	21000	66365290
Color	Steve Rash	13	98	15	281	Adam Arkin	691	8500000
Color	Ben Lewin	337	95	20	365	Adam Arkin	422	5997134
Color	Roland Emmerich	185	154	776	933	Adam Baldwin	10000	306124059
Color	Roland Emmerich	192	142	776	1000	Adam Baldwin	13000	113330342
Color	Chris Robinson	54	105	49	104	Adam Boyer	680	21160089
Color	Peter Jackson	645	182	0	773	Adam Brown	5000	303001229
Color	Peter Jackson	509	186	0	773	Adam Brown	5000	258355354
Color	Peter Jackson	422	164	0	773	Adam Brown	5000	255108370
Color	Michael McGowan	47	98	12	149	Adam Butcher	393	795126

Tabla de datos “moviedata.xlsx”
Fuente: “Escuela de Ingeniería UCN”

Las celdas vacías y erróneas contenían los siguientes N/A, #ERROR, #N/I, ?, #VALOR!, #NODATA. Lo que decidimos hacer en este caso fue:

- Pasar datos tipo texto no conocidos como “UNKNOWN”.
- Pasar datos tipo número no conocidos como “0”.
- “?,-1,0” no se eliminaron, dado que no afectaba a los datos el cambiarlos.

2. Limpieza de los datos

También añadimos funciones que nos sirven para la limpieza luego de leer los datos de cada columna, definiendo los valores NULL como 'nan' que son las siguientes:

- `cleanColor()`: Limpia los datos de la columna "color"
 - Reemplaza el string 'ERROR' por UNKNOWN
 - Reemplaza el cast 'nan' por UNKNOWN
 - Corta los espacios antes de que empiece el string.
- `cleanName()`: Limpia los datos de las columnas tipo texto.
 - Reemplaza el string '#NODATA' por NULL
 - Reemplaza el cast 'nan' por UNKNOWN
 - Reemplaza el string '?' por UNKNOWN
 - Corta los espacios antes de que empiece el string y transforma los errores traduciendo los a caracteres latinos en UTF-8.
- `cleanMovie()`: Limpia los datos de la columna "title_movie".
 - Elimina el último carácter 'Â' del string
- `cleanDigit()`: Limpia los datos de las columnas que contienen dígitos enteros.
 - Reemplaza el string '#NODATA' por UNKNOWN
 - Reemplaza el string 'nan' por UNKNOWN
 - Reemplaza el string 'N/I' por UNKNOWN
- `cleanDecimal()`: Limpia los datos de las columnas que contienen dígitos flotantes.
 - Reemplaza el string 'nan' por UNKNOWN

También hubo columnas que no se limpiaron, esto fué porque los datos están 100% funcionales y sin problemas para un posterior análisis y agregación a la base de datos, las columnas son las siguientes:

- `genres`
- `num_voted_users`
- `cast_total_facebook_likes`
- `movie_imdb_link`
- `language`

- country
- content_rating
- budget
- title_year
- imdb_score
- movie_facebook_likes

Además hubo columnas que fusionamos, en este caso “num_voted_users” y “cast_total_facebook_likes” para una película en particular tenía una variación muy pequeña de votos en aproximadamente 2 o 3 votos por película, por lo que deducimos que no eran datos altamente representativos y no generaban un cambio al momento de agregarla a la base de datos, esto nos ayuda a que haya un ahorro de cantidad de datos ingresados.

movie_title	num_voted_users	cast_total_facebook_likes
DisturbiaÂ	186984	2287
DisturbiaÂ	186982	2288
The WacknessÂ	27266	2748
The DictatorÂ	213863	1375
Elizabeth: The Golden AgeÂ	54787	16899
Sucker PunchÂ	197584	7067
RoboCopÂ	182899	14161
RoboCopÂ	182910	14160

Lo que decidimos en este caso fue obtener un promedio entre las filas con películas del mismo nombre y reducirlas a una sola fila.

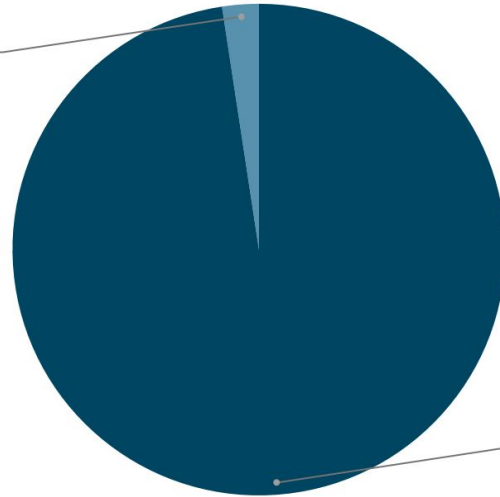
3. Comparación de datos

Luego de hacer toda la limpieza de datos anterior, decidimos comparar los datos que teníamos al principio, con los datos nuevos dentro del mismo dataframe modificado y llegamos a la siguiente conclusión:

Porcentaje de datos perdidos con respecto al original

Filtrado

2,4%



Original

97,6%

Cantidad de datos Originales vs Cantidad de datos filtrados

Fuente: Elaboración de Equipo

La cantidad de datos originales era de 5.043 columnas, mientras que luego de haber limpiado y filtrado todos los datos quedaron, disminuyó la cantidad de filas que teníamos al principio a 4.919. Esto significa que hubo un 2,4% de pérdida de información.

4. Exportación de archivos

```
movies.txt
http://www.imdb.com/title/tt1211956/?ref=fn_tt_tt_1,Escape Plan,Mikael Häfström,Color,286,115,101,25121291,85720,0,4230,0,1,cia
agent|escape|muslim|prison|ship,Action|Crime|Mystery|Sci-Fi|Thriller,279,English,USA,2,35,6,7,34000,2013,50000000,R,0,379,3376
http://www.imdb.com/title/tt1085492/?ref=fn_tt_tt_1,The Prince,Brian A Miller,Color,46,93,32,0,48458,0,9271,0,4,assassin|fight|
mechanic|rescue|rival,Action|Thriller,101,English,USA,2,35,4,6,0,2014,18000000,R,1,379,916
http://www.imdb.com/title/tt1798684/?ref=fn_tt_tt_1,Southpaw,Antoine Fuqua,Color,385,124,845,52418902,1418,0,19,0,0,boxer|boxing
training|death of wife|male in shower|page,Drama|Sport,277,English,USA,2,35,7,5,44000,2015,30000000,R,2,379,3549
http://www.imdb.com/title/tt0095863/?ref=fn_tt_tt_1,Phantasm II,Don Coscarelli,Color,81,97,152,7282851,67296,0,27842,0,1,cemetery|
female nudity|mortuary|sphere|tall man,Action|Fantasy|Horror|Sci-Fi|Thriller,100,English,USA,1,85,6,5,812,1988,3000000,R,3,2091,4113
http://www.imdb.com/title/tt0425151/?ref=fn_tt_tt_1,Jimmy and Judy,Randall Rubin,Color,23,99,0,0,2719,0,5900,0,2,police|revenge|sex|
suburb|video camera,Crime|Drama|Thriller,26,English,USA,0,0,6,2,138,2006,1000000,R,4,1108,4114
http://www.imdb.com/title/tt0403118/?ref=fn_tt_tt_1,Brave New Girl,Bobby Roth,Color,1,120,40,0,11003,0,22006,0,2,based on a book|
based on a novel|made for tv movie|music school|support,Drama|Family|Music,27,English,USA,0,0,5,0,47,2004,-1,PG-13,5,2092,4115
http://www.imdb.com/title/tt3319920/?ref=fn_tt_tt_1,Regression,Alejandro Amenábar,Color,140,106,448,54734,2843,0,102,0,1,inspired by
true events|memory|Minnesota|police|year 1990,Crime|Drama|Mystery|Thriller,62,English,Spain,2,35,5,7,0,2015,-1,R,6,2092,1316
http://www.imdb.com/title/tt0796375/?ref=fn_tt_tt_1,You Kill Me,John Dahl,Color,121,93,131,2426851,51842,0,9330,0,1,alcoholic|
buffalo new york|claim in title|embalming|mortuary,Comedy|Crime|Romance|
Thriller,76,English,USA,2,55,6,5,830,2007,4000000,R,7,2093,4116
http://www.imdb.com/title/tt4591310/?ref=fn_tt_tt_1,Xi you ji zhi: Sun Wukong san da Baigu Jing,Pou-Soi
Cheang,Color,14,119,3,0,14786,0,16768,0,1,buddhism|demon|journey to the west|monk|monkey king,Action|Adventure|
Fantasy,9,English,China,2,35,6,0,426,2016,68005000,I,8,2094,4117
http://www.imdb.com/title/tt0271219/?ref=fn_tt_tt_1,Tadpole,Gary Winick,Color,91,78,56,2882062,9638,0,39,0,2,best friend|boy|french|
friend|love,Comedy|Drama|Romance,101,English,USA,1,85,6,3,132,2000,150000,PG-13,9,2095,3754
http://www.imdb.com/title/tt0426931/?ref=fn_tt_tt_1,August Rush,Kirsten Sheridan,Color,148,114,59,31655091,87351,0,43917,0,0,baby|
cellist|genius|prodigy|rhapsody,Drama|Music,364,English,USA,2,35,7,5,18000,2007,30000000,PG,10,2096,3332
http://www.imdb.com/title/tt0486822/?ref=fn_tt_tt_1,Disturbia,D.J. Caruso,Color,253,105,154,80050171,48346,0,31014,0,0,binoculars|
electronic tag|house arrest|neighbor|watching someone,Drama|Mystery|
Thriller,491,English,USA,1,85,6,9,0,2007,20000000,PG-13,11,2097,4118
http://www.imdb.com/title/tt1082886/?ref=fn_tt_tt_1,The Wackness,Jonathan Levine,Color,147,99,129,2077046,212085,0,16034,0,2,ice
cream|marijuana|new york city|summer|therapy,Comedy|Drama|Romance,75,English,USA,2,35,7,0,0,2008,6000000,R,12,2097,4119
http://www.imdb.com/title/tt1645170/?ref=fn_tt_tt_1,The Dictator,Larry Charles,Color,313,99,119,59617068,7277,0,417,0,0,dictator|
```

```
actors.txt
0,Sylvester Stallone,13000
1,Bruce Willis,13000
2,Jake Gyllenhaal,15000
3,Angus Scrimm,674
4,Nicole Randall Johnson,363
5,Virginia Madsen,912
6,Emma Watson,9000
7,Philip Baker Hall,497
8,Li Gong,879
9,Bebe Neuwirth,376
10,Robin Williams,49000
11,Sarah Roemer,884
12,Mary-Kate Olsen,976
13,Sayed Badreya,600
14,Eddie Redmayne,13000
15,Jon Hamm,4000
16,Gary Oldman,10000
17,Vidya Balan,464
18,Gerard Butler,18000
19,Tom Hardy,27000
20,Tina Fey,2000
21,Brendan Fraser,3000
22,Michael Raymond-James,788
23,Jack E. Currenton,847
24,Robert Downey Jr.,21000
25,Mako,691
26,W. Earl Brown,422
27,Will Smith,10000
28,Heath Ledger,13000
29,T.I.,600
```

Archivos finales del proyecto

Fuente: Elaboración propia

Aquí se tienen todos los archivos que se generan luego de ejecutar el código principal. Se generan los archivos “movies.txt” y “actors.txt” para su posterior carga a la base de datos. Lo hicimos de esta manera (Importación-Exportación) para evitar la espera que se procesen los datos mientras hacemos pruebas al programar. Esto puede mejorar haciendo uso de hilos.

5. Carga a la bases de datos

Finalmente para la elección de tablas decidimos crear 2 tablas mostradas a continuación:

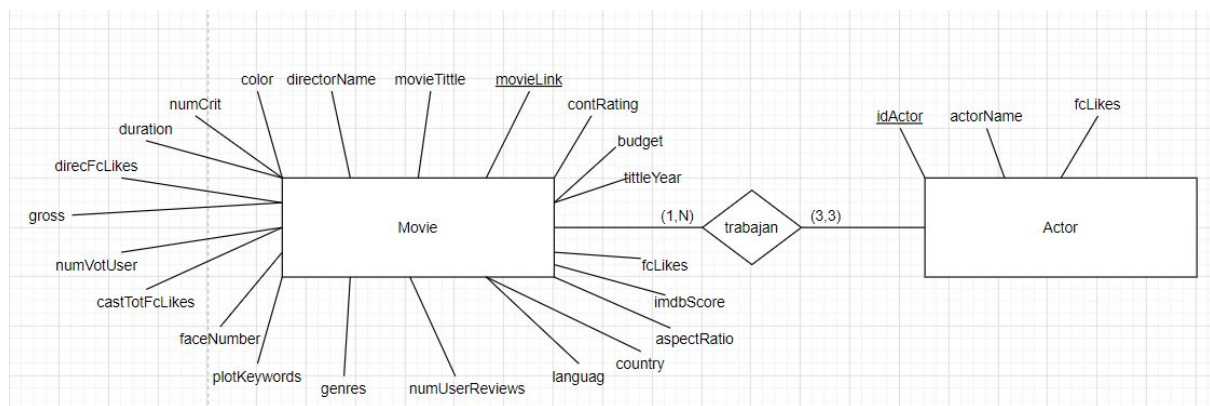


Diagrama Entidad-Relacionamiento

Fuente: Elaboración de Equipo