# DALHOUSIE UNIVERSITY

## CSCI 5408 – Data Management, Warehousing, Analytics

## Assignment 3

**Work done by,**

**Name: Guturu Rama Mohan Vishnu**

**Banner ID: B00871849**

**Email: rm286720@dal.ca**

# DECLARATION

**I, Guturu Rama Mohan Vishnu, declare that in assignment 3 of CSCI 5408 course, data scrapping is not done programmatically or using any online or offline tools. However, the webpages or the domain mentioned in this document are visited manually, and some useful information is gathered for education purpose only. Information, such as email, personal contact numbers, or names of people are not extracted. The course instructor or the Faculty of Computer Science cannot be held responsible for any misuse of the extracted data.**

**Problem #2:** MySQL and Java – Implementation of 2 phase locking

**Solution:**

The datasets I found in the given link are:

   (a) AnnualTicketSales
   (b) HighestGrossers
   (c) PopularCreativeTypes
   (d) TopDistributors
   (e) TopGenres
   (f) TopGrossingRatings
   (g) TopGrossingSources
   (h) TopProductionMethods
   (i) WideReleasesCount

## Cleaning and Transformation

## AnnualTicketSales

There is no primary key for this table.

(1) There are also no blank cells or blank columns in this dataset.

(2) But the columns contain values which have "$" sign in them before the actual value. In order to load the data into database and in order to compare the values, we have to remove the $ sign.



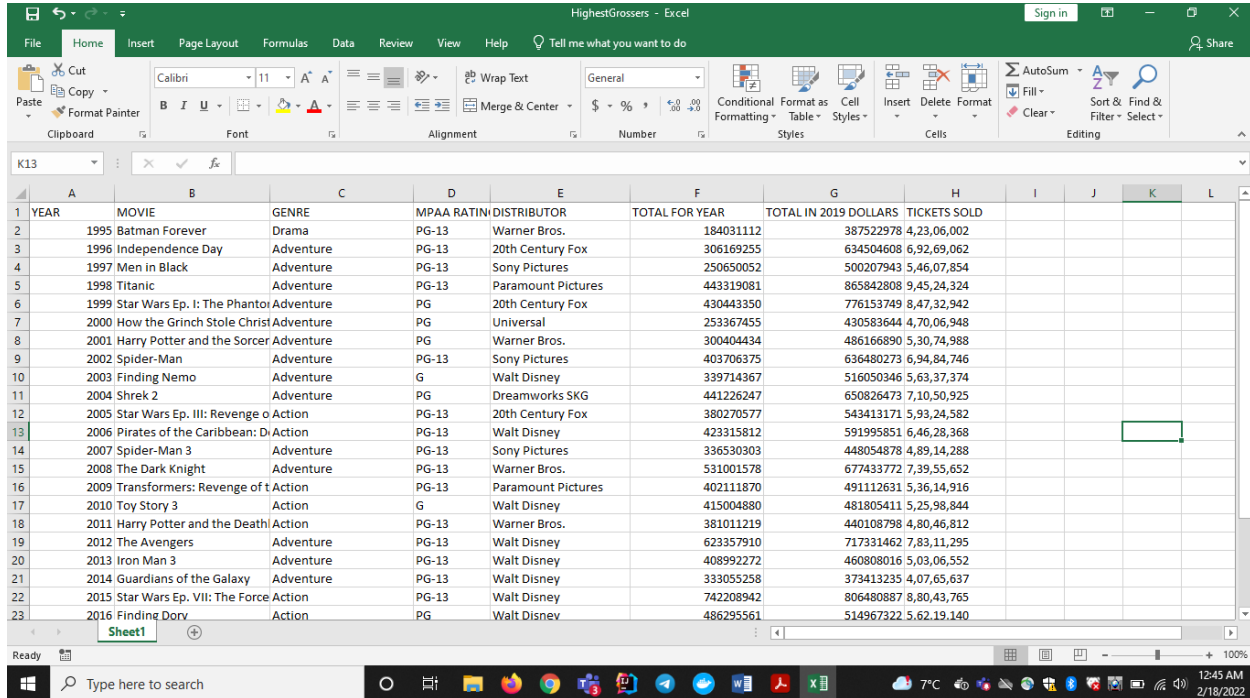(3) So, I am replacing it with empty string "", i.e., I am removing that character from the field.

## HighestGrossers

The Primary Key for this table is "MOVIE" column.

(1) We have few blank fields in this data set. We replace them with NaN since they are in text type column.



| | MOVIE | | | | | | |
|---|---|---|---|---|---|---|---|
| 2004 | Shrek 2 | Adventure | PG | Dreamworks SKG | $441,226,247 | $650,826,473 | 7,10,50,925 |
| 2005 | Star Wars Ep. III: Revenge o | Action | PG-13 | 20th Century Fox | $380,270,577 | $543,413,171 | 5,93,24,582 |
| 2006 | Pirates of the Caribbean: D | Action | PG-13 | Walt Disney | $423,315,812 | $591,995,851 | 6,46,28,368 |
| 2007 | Spider-Man 3 | Adventure | PG-13 | Sony Pictures | $336,530,303 | $448,054,878 | 4,89,14,288 |
| 2008 | The Dark Knight | Adventure | PG-13 | Warner Bros. | $531,001,578 | $677,433,772 | 7,39,55,652 |
| 2009 | Transformers: Revenge of t | Action | PG-13 | Paramount Pictures | $402,111,870 | $491,112,631 | 5,36,14,916 |
| 2010 | Toy Story 3 | Action | G | Walt Disney | $415,004,880 | $481,805,411 | 5,25,98,844 |
| 2011 | Harry Potter and the Death | Action | PG-13 | Warner Bros. . | $381,011,219 | $440,108,798 | 4,80,46,812 |
| 2012 | The Avengers | Adventure | PG-13 | Walt Disney | $623,357,910 | $717,331,462 | 7,83,11,295 |
| 2013 | Iron Man 3 | Adventure | PG-13 | Walt Disney | $408,992,272 | $460,808,016 | 5,03,06,552 |
| 2014 | Guardians of the Galaxy | Adventure | PG-13 | Walt Disney | $333,055,258 | $373,413,235 | 4,07,65,637 |
| 2015 | Star Wars Ep. VII: The Force | Action | PG-13 | Walt Disney | $742,208,942 | $806,480,887 | 8,80,43,765 |
| 2016 | Finding Dory | Action | PG | Walt Disney | $486,295,561 | $514,967,322 | 5,62,19,140 |
| 2017 | Star Wars Ep. VIII: The Last | Action | PG-13 | Walt Disney | $517,218,368 | $528,173,936 | 5,76,60,910 |
| 2018 | Black Panther | Action | PG-13 | Walt Disney | $700,059,566 | $703,901,821 | 7,68,45,177 |
| 2019 | Avengers: Endgame | NaN | PG-13 | Walt Disney | $858,373,000 | $858,373,002 | 9,37,08,843 |
| 2020 | Bad Boys For Life | NaN | R | Sony Pictures | $204,417,855 | $204,417,848 | 2,23,16,359 |
| 2021 | Shang-Chi and the Legend o | NaN | PG-13 | Walt Disney | $224,226,704 | $224,226,704 | 2,44,78,897 |

(2) As saw in the first table, few columns contain values which have "$" sign in them before the actual value. So, we do the same thing which we did before, in the process of data cleaning.



| YEAR | MOVIE | GENRE | MPAA RATING | DISTRIBUTOR | TOTAL FOR YEAR | TOTAL IN 2019 DOLLARS | TICKETS SOLD |
|---|---|---|---|---|---|---|---|
| 1995 | Batman Forever | Drama | PG-13 | Warner Bros. | 184031112 | 387522978 | 4,23,06,002 |
| 1996 | Independence Day | Adventure | PG-13 | 20th Century Fox | 306169255 | 634504608 | 6,92,69,062 |
| 1997 | Men in Black | Adventure | PG-13 | Sony Pictures | 250650052 | 500207943 | 5,46,07,854 |
| 1998 | Titanic | Adventure | PG-13 | Paramount Pictures | 443319081 | 865842808 | 9,45,24,324 |
| 1999 | Star Wars Ep. I: The Phantom | Adventure | PG | 20th Century Fox | 430443350 | 776153749 | 8,47,32,942 |
| 2000 | How the Grinch Stole Christm | Adventure | PG | Universal | 253367455 | 430583644 | 4,70,06,948 |
| 2001 | Harry Potter and the Sorcer | Adventure | PG | Warner Bros. | 300404434 | 486166890 | 5,30,74,988 |
| 2002 | Spider-Man | Adventure | PG-13 | Sony Pictures | 403706375 | 636480273 | 6,94,84,746 |
| 2003 | Finding Nemo | Adventure | G | Walt Disney | 339714367 | 516050346 | 5,63,37,374 |
| 2004 | Shrek 2 | Adventure | PG | Dreamworks SKG | 441226247 | 650826473 | 7,10,50,925 |
| 2005 | Star Wars Ep. III: Revenge o | Action | PG-13 | 20th Century Fox | 380270577 | 543413171 | 5,93,24,582 |
| 2006 | Pirates of the Caribbean: D | Action | PG-13 | Walt Disney | 423315812 | 591995851 | 6,46,28,368 |
| 2007 | Spider-Man 3 | Adventure | PG-13 | Sony Pictures | 336530303 | 448054878 | 4,89,14,288 |
| 2008 | The Dark Knight | Adventure | PG-13 | Warner Bros. | 531001578 | 677433772 | 7,39,55,652 |
| 2009 | Transformers: Revenge of t | Action | PG-13 | Paramount Pictures | 402111870 | 491112631 | 5,36,14,916 |
| 2010 | Toy Story 3 | Action | G | Walt Disney | 415004880 | 481805411 | 5,25,98,844 |
| 2011 | Harry Potter and the Deathl | Action | PG-13 | Warner Bros. | 381011219 | 440108798 | 4,80,46,812 |
| 2012 | The Avengers | Adventure | PG-13 | Walt Disney | 623357910 | 717331462 | 7,83,11,295 |
| 2013 | Iron Man 3 | Adventure | PG-13 | Walt Disney | 408992272 | 460808016 | 5,03,06,552 |
| 2014 | Guardians of the Galaxy | Adventure | PG-13 | Walt Disney | 333055258 | 373413235 | 4,07,65,637 |
| 2015 | Star Wars Ep. VII: The Force | Action | PG-13 | Walt Disney | 742208942 | 806480887 | 8,80,43,765 |
| 2016 | Finding Dory | Action | PG | Walt Disney | 486295561 | 514967322 | 5,62,19,140 |

(3) The column "TICKETS SOLD" has the integer values but as form of text. So, in order for the database to detect it as a number, I have changed all the fields to number format.

| YEAR | MOVIE | GENRE | MPAA RATING | DISTRIBUTOR | TOTAL FOR YEAR | TOTAL IN 2019 DOLLARS | TICKETS SOLD |
|---|---|---|---|---|---|---|---|
| 1995 | Batman Forever | Drama | PG-13 | Warner Bros. | 184031112 | 387522978 | 42306002 |
| 1996 | Independence Day | Adventure | PG-13 | 20th Century Fox | 306169255 | 634504608 | 69269062 |
| 1997 | Men in Black | Adventure | PG-13 | Sony Pictures | 250650052 | 500207943 | 54607854 |
| 1998 | Titanic | Adventure | PG-13 | Paramount Pictures | 443319081 | 865842808 | 94524324 |
| 1999 | Star Wars Ep. I: The Phantor | Adventure | PG | 20th Century Fox | 430443350 | 776153749 | 84732942 |
| 2000 | How the Grinch Stole Christ | Adventure | PG | Universal | 253367455 | 430583644 | 47006948 |
| 2001 | Harry Potter and the Sorcer | Adventure | PG | Warner Bros. | 300404434 | 486166890 | 53074988 |
| 2002 | Spider-Man | Adventure | PG-13 | Sony Pictures | 403706375 | 636480273 | 69484746 |
| 2003 | Finding Nemo | Adventure | G | Walt Disney | 339714367 | 516050346 | 56337374 |
| 2004 | Shrek 2 | Adventure | PG | Dreamworks SKG | 441226247 | 650826473 | 71050925 |
| 2005 | Star Wars Ep. III: Revenge o | Action | PG-13 | 20th Century Fox | 380270577 | 543413171 | 59324582 |
| 2006 | Pirates of the Caribbean: D | Action | PG-13 | Walt Disney | 423315812 | 591995851 | 64628368 |
| 2007 | Spider-Man 3 | Adventure | PG-13 | Sony Pictures | 336530303 | 448054878 | 48914288 |
| 2008 | The Dark Knight | Adventure | PG-13 | Warner Bros. | 531001578 | 677433772 | 73955652 |
| 2009 | Transformers: Revenge of t | Action | PG-13 | Paramount Pictures | 402111870 | 491112631 | 53614916 |
| 2010 | Toy Story 3 | Action | G | Walt Disney | 415004880 | 481805411 | 52598844 |
| 2011 | Harry Potter and the Deathl | Action | PG-13 | Warner Bros. | 381011219 | 440108798 | 48046812 |
| 2012 | The Avengers | Adventure | PG-13 | Walt Disney | 623357910 | 717331462 | 78311295 |
| 2013 | Iron Man 3 | Adventure | PG-13 | Walt Disney | 408992272 | 460808016 | 50306552 |
| 2014 | Guardians of the Galaxy | Adventure | PG-13 | Walt Disney | 333055258 | 373413235 | 40765637 |
| 2015 | Star Wars Ep. VII: The Force | Action | PG-13 | Walt Disney | 742208942 | 806480887 | 88043765 |
| 2016 | Finding Dory | Action | PG | Walt Disney | 486295561 | 514967322 | 56219140 |

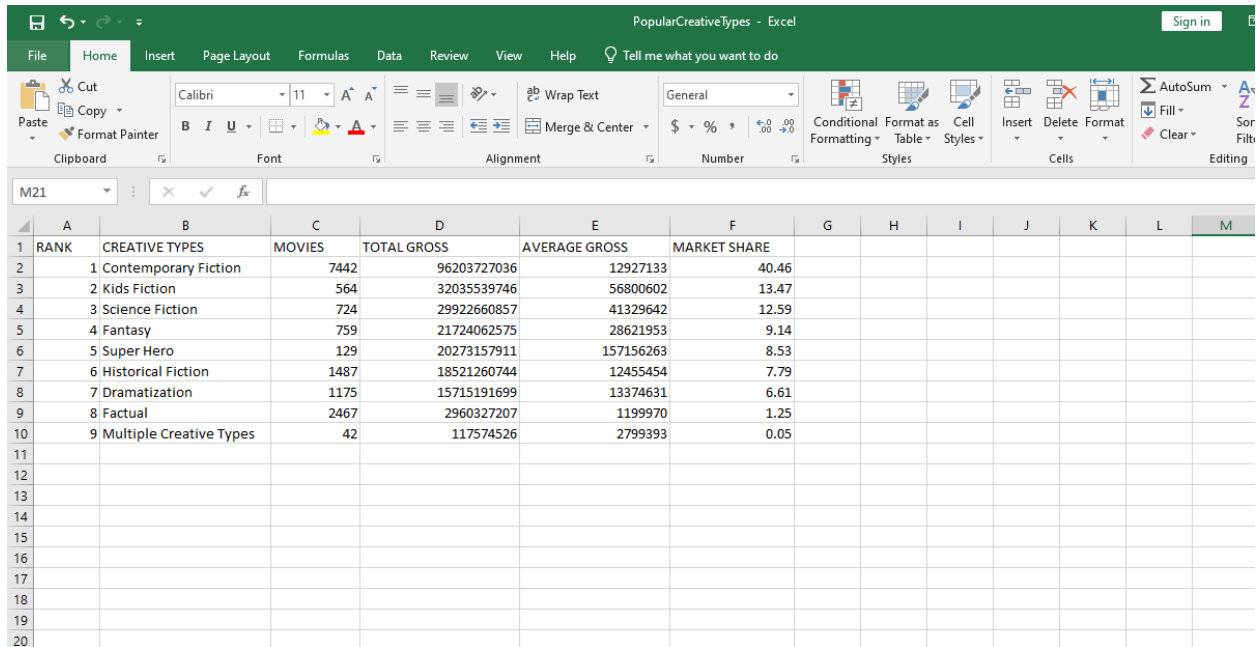(4) Also, there are no duplicate values in the primary key column.

## PopularCreativeTypes

The Primary Key for this table is "CREATIVE TYPES".

(1) There aren't any blank cells in this data set.

(2) Few columns contain values which have "$" sign in them before the actual value. So, we do the same thing which we did before, in the process of data cleaning. I removed the signs.

| RANK | CREATIVE TYPES | MOVIES | TOTAL GROSS | AVERAGE GROSS | MARKET SHARE |
|------|----------------|--------|-------------|---------------|--------------|
| 1 | Contemporary Fiction | 7,442 | 96203727036 | 12927133 | 40.46% |
| 2 | Kids Fiction | 564 | 32035539746 | 56800602 | 13.47% |
| 3 | Science Fiction | 724 | 29922660857 | 41329642 | 12.59% |
| 4 | Fantasy | 759 | 21724062575 | 28621953 | 9.14% |
| 5 | Super Hero | 129 | 20273157911 | 157156263 | 8.53% |
| 6 | Historical Fiction | 1,487 | 18521260744 | 12455454 | 7.79% |
| 7 | Dramatization | 1,175 | 15715191699 | 13374631 | 6.61% |
| 8 | Factual | 2,467 | 2960327207 | 1199970 | 1.25% |
| 9 | Multiple Creative Types | 42 | 117574526 | 2799393 | 0.05% |

(3) Also, there is a column which has a symbol % in it. So, we remove it too for consistency.



| RANK | CREATIVE TYPES | MOVIES | TOTAL GROSS | AVERAGE GROSS | MARKET SHARE |
|------|----------------|--------|-------------|---------------|--------------|
| 1 | Contemporary Fiction | 7,442 | 96203727036 | 12927133 | 40.46 |
| 2 | Kids Fiction | 564 | 32035539746 | 56800602 | 13.47 |
| 3 | Science Fiction | 724 | 29922660857 | 41329642 | 12.59 |
| 4 | Fantasy | 759 | 21724062575 | 28621953 | 9.14 |
| 5 | Super Hero | 129 | 20273157911 | 157156263 | 8.53 |
| 6 | Historical Fiction | 1,487 | 18521260744 | 12455454 | 7.79 |
| 7 | Dramatization | 1,175 | 15715191699 | 13374631 | 6.61 |
| 8 | Factual | 2,467 | 2960327207 | 1199970 | 1.25 |
| 9 | Multiple Creative Types | 42 | 117574526 | 2799393 | 0.05 |

(4) Movies column has comma in the values. To maintain consistency, I removed them too.



(5) There are no duplicate values in the primary key column.

## TopDistributors

The Primary Key for this table is "DISTRIBUTORS".

(1) There aren't any blank cells in this data set.

(2) Few columns contain values which have "$" sign in them before the actual value. So, we do the same thing which we did before, in the process of data cleaning. I removed the signs.

(3) Also, there is a column which has a symbol % in it. So, we remove it too for consistency.



(4) There are no duplicate values in the primary key column.

(5) There are a few empty columns in this dataset which can be seen while importing data in to the database. I deleted those columns too.

**TopGenres**

The Primary Key for this table is "GENRES".

(1) There aren't any blank cells in this data set.

(2) Few columns contain values which have "$" sign in them before the actual value. So, we do the same thing which we did before, in the process of data cleaning. I removed the signs.



| RANK | GENRES | MOVIES | TOTAL GROSS | AVERAGE GROSS | MARKET SHARE |
|---|---|---|---|---|---|
| 1 | Adventure | 1,102 | 64529536530 | 58556748 | 27.14% |
| 2 | Action | 1,098 | 49339974493 | 44936224 | 20.75% |
| 3 | Drama | 5,479 | 35586177269 | 6495013 | 14.97% |
| 4 | Comedy | 2,418 | 33687992318 | 13932172 | 14.17% |
| 5 | Thriller/Suspense | 1,186 | 19810201102 | 16703374 | 8.33% |
| 6 | Horror | 716 | 13430378699 | 18757512 | 5.65% |
| 7 | Romantic Comedy | 630 | 10480124374 | 16635118 | 4.41% |
| 8 | Musical | 201 | 4293988317 | 21363126 | 1.81% |
| 9 | Documentary | 2,415 | 2519513142 | 1043277 | 1.06% |
| 10 | Black Comedy | 213 | 2185433323 | 10260250 | 0.92% |

(3) Also, there is a column which has a symbol % in it. So, we remove it too for consistency.

(4) Movies column has comma in the values. To maintain consistency, I removed them too.



(5) There are no duplicate values in the primary key column.

**TopGrossingRatings**

The Primary Key for this table is "MPAA RATINGS".

(1) There aren't any blank cells in this data set.

(2) Few columns contain values which have "$" sign in them before the actual value. So, we do the same thing which we did before, in the process of data cleaning. I removed the signs.



| RANK | MPAA RATINGS | MOVIES | TOTAL GROSS | AVERAGE GROSS | MARKET SHARE |
|------|--------------|--------|-------------|---------------|--------------|
| 1 | PG-13 | 3,243 | 1.13525E+11 | 35006102 | 47.75% |
| 2 | R | 5,480 | 63497164978 | 11587074 | 26.71% |
| 3 | PG | 1,535 | 49124317794 | 32002813 | 20.66% |
| 4 | G | 395 | 9572240391 | 24233520 | 4.03% |
| 5 | Not Rated | 5,820 | 1918358283 | 329615 | 0.81% |
| 6 | NC-17 | 24 | 44850139 | 1868756 | 0.02% |
| 7 | Open | 5 | 5489687 | 1097937 | 0.00% |
| 8 | GP | 7 | 552618 | 78945 | 0.00% |

(3) Also, there is a column which has a symbol % in it. So, we remove it too for consistency.

(4) Movies column has comma in the values. To maintain consistency, I removed them too.



(5) There are no duplicate values in the primary key column.

**TopGrossingSources**

The Primary Key for this table is "SOURCES".

(1) There aren't any blank cells in this data set.

(2) Few columns contain values which have "$" sign in them before the actual value. So, we do the same thing which we did before, in the process of data cleaning. I removed the signs.



| RANK | SOURCES | MOVIES | TOTAL GROSS | AVERAGE GROSS | MARKET SHARE |
|---|---|---|---|---|---|
| 1 | Original Screenplay | 7,946 | 1.06375E+11 | 13387264 | 44.74% |
| 2 | Based on Fiction Book/Short Story | 2,150 | 47005613207 | 21863076 | 19.77% |
| 3 | Based on Comic/Graphic Novel | 249 | 23369989130 | 93855378 | 9.83% |
| 4 | Remake | 328 | 12832659970 | 39123963 | 5.40% |
| 5 | Based on Real Life Events | 3,225 | 11398356297 | 3534374 | 4.79% |
| 6 | Based on TV | 231 | 11305006312 | 48939421 | 4.75% |
| 7 | Based on Factual Book/Article | 295 | 7443681990 | 25232820 | 3.13% |
| 8 | Spin-Off | 41 | 3833128331 | 93490935 | 1.61% |
| 9 | Based on Folk Tale/Legend/Fairytale | 78 | 3406118495 | 43668186 | 1.43% |
| 10 | Based on Play | 271 | 2111190923 | 7790372 | 0.89% |

(3) Also, there is a column which has a symbol % in it. So, we remove it too for consistency.

(4) Movies column has comma in the values. To maintain consistency, I removed them too.



(5) There are no duplicate values in the primary key column.

# TopProductionMethods

The Primary key for this table is "PRODUCTION METHODS".

(1) There aren't any blank cells in this data set.

(2) Few columns contain values which have "$" sign in them before the actual value. So, we do the same thing which we did before, in the process of data cleaning. I removed the signs.



| RANK | PRODUCTION METHODS | MOVIES | TOTAL GROSS | AVERAGE GROSS | MARKET SHARE |
|---|---|---|---|---|---|
| 1 | Live Action | 14,613 | 1.79637E+11 | 12292972 | 75.56% |
| 2 | Animation/Live Action | 264 | 30346622254 | 114949327 | 12.76% |
| 3 | Digital Animation | 365 | 23920180508 | 65534741 | 10.06% |
| 4 | Hand Animation | 164 | 2960497487 | 18051814 | 1.25% |
| 5 | Stop-Motion Animation | 37 | 676490120 | 18283517 | 0.28% |
| 6 | Multiple Production Methods | 26 | 43728300 | 1681858 | 0.02% |
| 7 | Rotoscoping | 4 | 8468385 | 2117096 | 0.00% |

(3) Also, there is a column which has a symbol % in it. So, we remove it too for consistency.

(4) Movies column has comma in the values. To maintain consistency, I removed them too.



(5) There are no duplicate values in the primary key column.

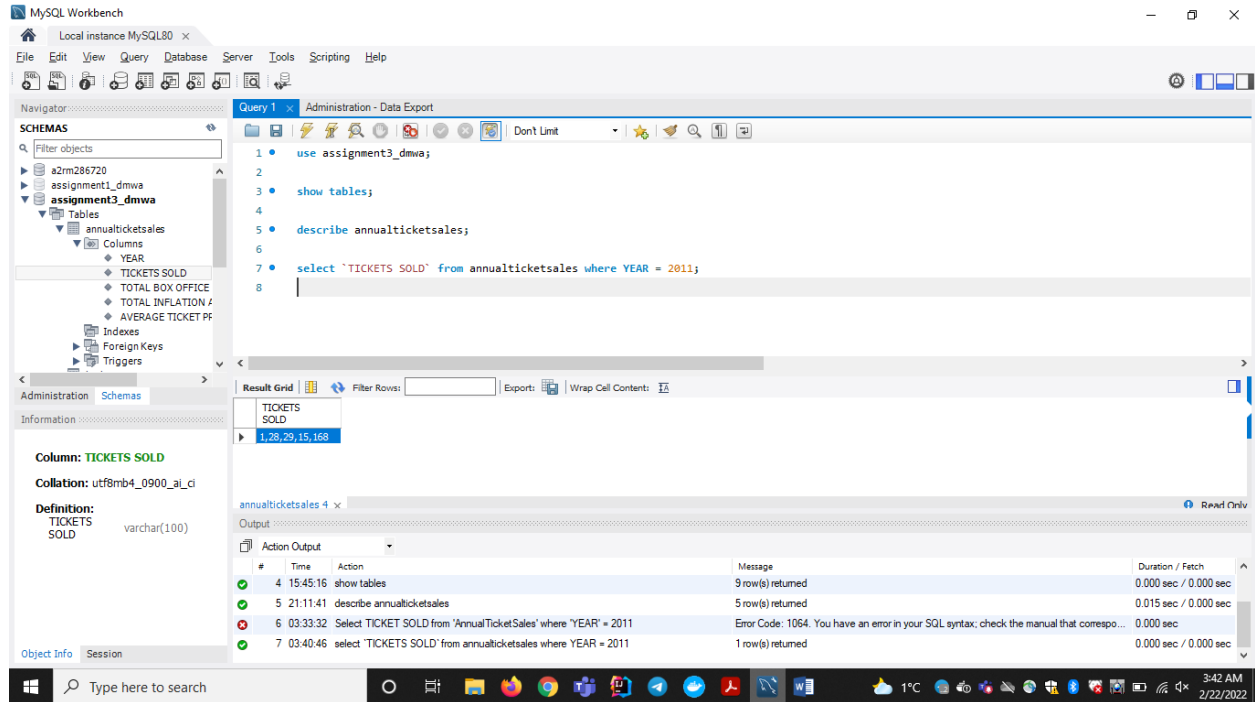**WideReleasesCount**

There is no primary key for this table.

(1) There are no blank values in this table.

(2) There are no special characters like **$** or **,** or **%**. Hence there is no cleaning to do in this dataset.
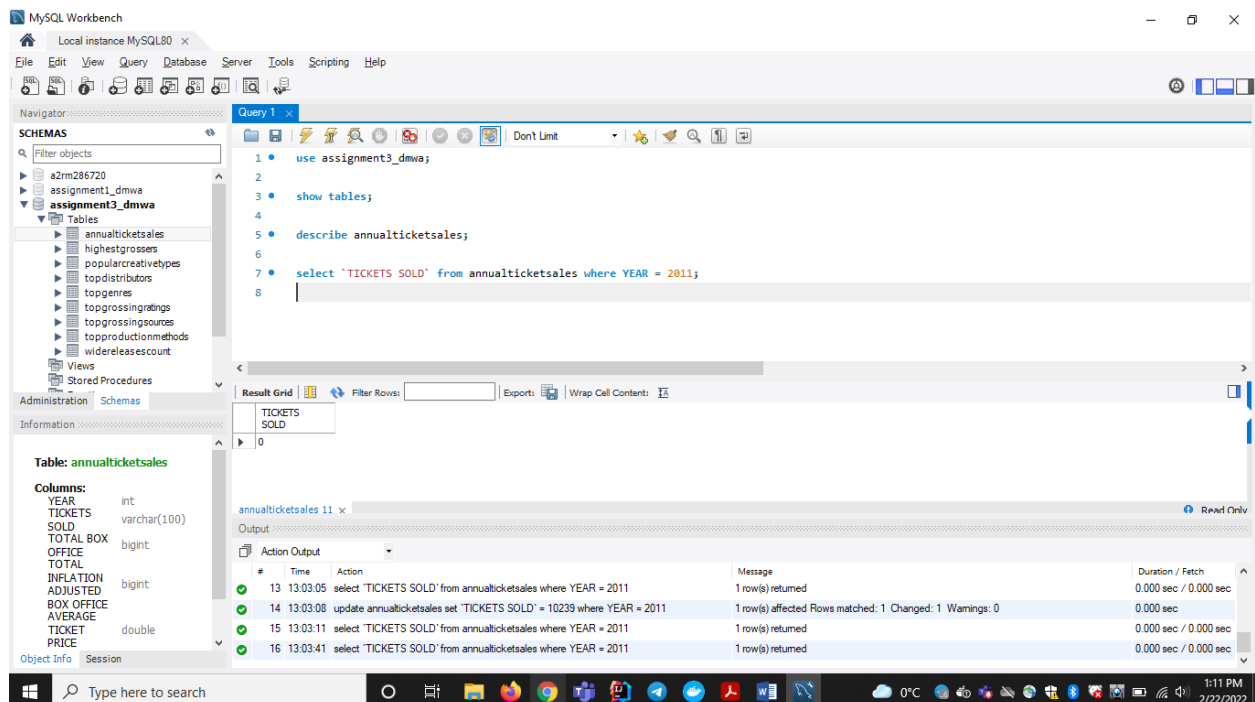
# Two Phase Locking Protocol

## Before Transaction:



## After Transaction:

# Proof that code got executed:



**Git URL for code:** https://git.cs.dal.ca/rguturu/csci-5408-w2022-b00871849-gutururamamohanvishnu/-/tree/main/assignment-3

**References:**

[1]     Tanwar, V., "*Two Phase Locking Protocol – GeeksforGeeks*," GeeksforGeeks [Online]. Available at: https://www.geeksforgeeks.org/two-phase-locking-protocol/ [Accessed: February 22, 2022].

[2]     Purswani, A., "*Difference between Shared Lock and Exclusive Lock – GeeksforGeeks*," GeeksforGeeks [Online]. Available at: https://www.geeksforgeeks.org/difference-between-shared-lock-and-exclusive-lock/ [Accessed: February 22, 2022].

[3]     "*ReentrantReadWriteLock - why can't reader acquire writer's lock?*," Stack Overflow [Online]. Available at: https://stackoverflow.com/questions/43099144/reentrantreadwritelock-why-cant-reader-acquire-writers-lock [Accessed: February 22, 2022].