# BIG DATA CAPSTONE

1.



2.

```
1 select * from moviesdata limit 10;
```

SQL

TEZ

Execute   Explain   Save as...

New Worksheet

**Query Process Results (Status: SUCCEEDED)**

Save results... ▾

Logs   Results

Filter columns...

previous   next

| moviesdata.name | moviesdata.rating | moviesdata.genre | moviesdata.year | moviesdata.released | moviesdata.sco |
|---|---|---|---|---|---|
| The Shining | R | Drama | 1980 | June 13, 1980 (United States) | 8.4 |
| The Blue Lagoon | R | Adventure | 1980 | July 2, 1980 (United States) | 5.8 |
| Star Wars: Episode V - The Empire Strikes Back | PG | Action | 1980 | June 20, 1980 (United States) | 8.7 |

### Search tables...

**Databases**

banala
banalas
default
  moviesdata

| name | STRING |
|---|---|
| rating | STRING |
| genre | STRING |
| year | INT |
| released | STRING |
| score | DOUBLE |
| votes | INT |
| director | STRING |
| writer | STRING |
| star | STRING |
| country | STRING |
| budget | INT |
| gross | INT |
| company | STRING |
| runtime | INT |

sample 07
sample 08

3.

Worksheet ✖   capstone 5 ✖   capstone 3 ✖

```
1 select rating,count(*) AS TOTALMOVIES from moviesdata group by rating order by TOTALMOVIES desc;
```

SQL

TEZ

Execute   Explain   Save as...

New Worksheet

100%

**Query Process Results (Status: SUCCEEDED)**

Save results... ▾

Logs   Results

Filter columns...

previous   next

| rating | totalmovies |
|---|---|
| R | 3697 |
| PG-13 | 2112 |
| PG | 1252 |
| Not Rated | 283 |

Search tables...

**Databases**

banala
banalas
default
foodmart
shyni
xademo

➢ From the above visualisation I can conclude that top 3 ratings maximum number of movies are

1. R (3697)
2. PG_13(2112)
3. PG(1252)

➢ That means maximum number of people from each rating who watched top 3 movies are r,pg_13,pg

➢ From this top most watched movies are 3697 .

4.

```
1 select genre,count(*) AS TOTALMOVIES from moviesdata group by genre order by TOTALMOVIES DESC limit 5
```

Search tables...

Databases
- banala
- banalas
- default
- foodmart
- shyni
- xademo

SQL

TEZ

Execute    Explain    Save as...    New Worksheet

**Query Process Results (Status: SUCCEEDED)**    Save results... ▾

Logs    Results
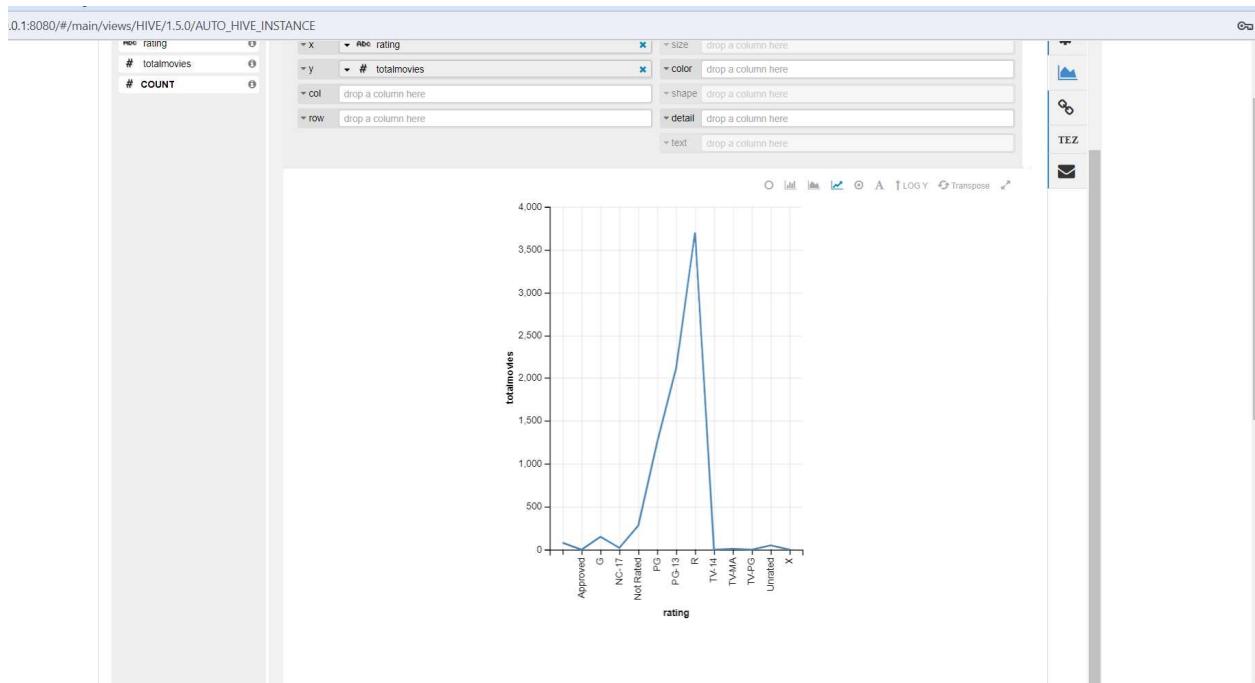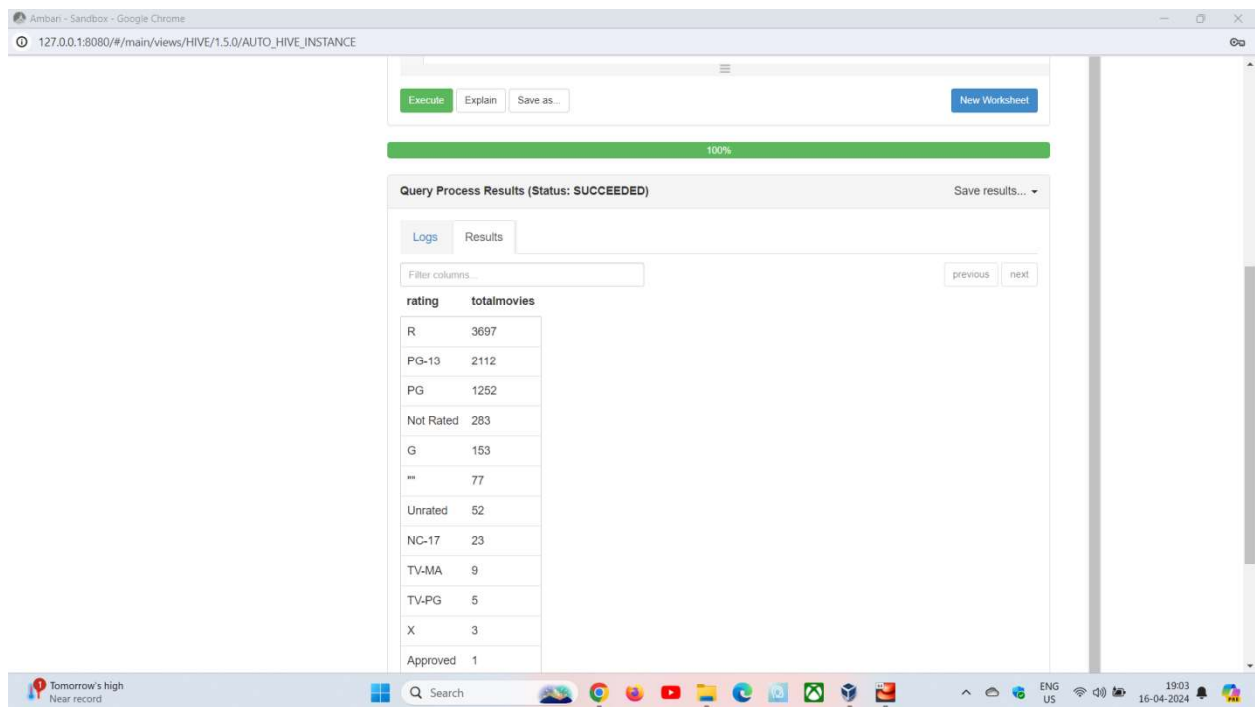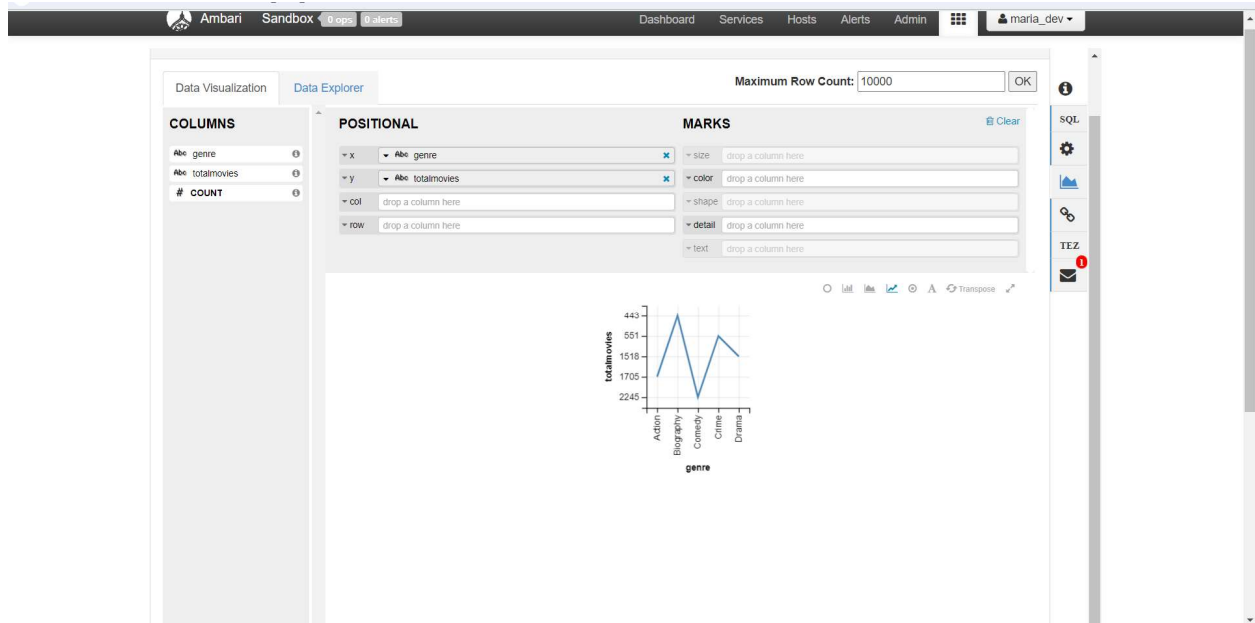
Filter columns...    previous    next

| genre | totalmovies |
| --- | --- |
| Comedy | 2245 |
| Action | 1705 |
| Drama | 1518 |
| Crime | 551 |
| Biography | 443 |

> From the above visualization I can conclude that from each genre top 3 movies are
>> 1.Comedy (2245)
>> 2.Action  (1705)
>> 3.Drama  (1518)
> From this most number of people watched the comedy movies .
> Most number of directors liked to make comedy movies.

5.

Ambari    Sandbox  0 ops  0 alerts          Dashboard    Services    Hosts    Alerts    Admin          maria_dev ▾

Hive    Query    Saved Queries    History    UDFs    Upload Table

**Database Explorer**    ⟳

default    ▾

Search tables...

Databases

🗄banala
🗄banalas
🗄default
🗄foodmart
🗄shyni
🗄xademo

**Query Editor**    ⤢

Worksheet ✖    capstone 3 ✖    capstone 4 ✖    capstone 13 ✖    capstone 5 ✖    capstone 12 ✖

1  select  count(year) AS TOP10YEARS,year  from moviesdata group by year order by TOP10YEARS desc limit

SQL
⚙
📊
🔗
TEZ
✉ 7

Execute    Explain    Save as...                    New Worksheet

100%

**Query Process Results (Status: SUCCEEDED)**    Save results... ▾

Logs    Results

---

Execute    Explain    Save as...                    New Worksheet

100%

**Query Process Results (Status: SUCCEEDED)**    Save results... ▾

Logs    Results

Filter columns...                    previous    next

| top10years | year |
|---|---|
| 200 | 1985 |
| 200 | 2019 |
| 200 | 1994 |
| 200 | 1987 |
| 200 | 1988 |
| 200 | 1991 |
| 200 | 1986 |
| 200 | 2018 |
| 200 | 1992 |
| 200 | 1989 |

- From the above visualisation I can conclude that  in the year of 1985  200 movies were released .
- In the years  of 2019,1994 ,1987,1988,1991,1986, 2018,1992,1989  total 200 movies were released .
- For top 10 years all the movies released were same .

6.

```
1 select name,score from moviesdata order by score desc;
```

Databases

banala
banalas
default
foodmart
shynii
xademo

TEZ

Execute    Explain    Save as...                                    New Worksheet

100%

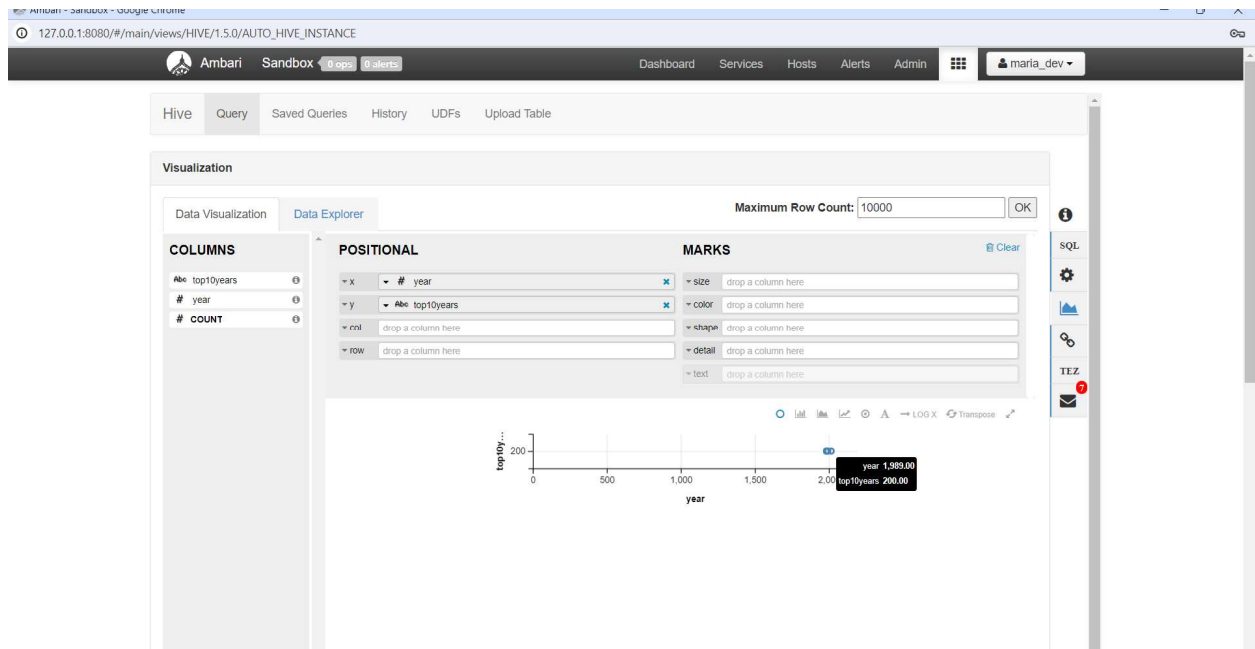**Query Process Results (Status: SUCCEEDED)**                        Save results... ▾

Logs    Results

Filter columns...                                          previous    next

| name | score |
|------|-------|
| The Shawshank Redemption | 9.3 |
| The Dark Knight | 9.0 |
| The Lord of the Rings: The Return of the King | 8.9 |
| Pulp Fiction | 8.9 |

7.

```
1  select name,votes from moviesdata order by votes desc limit 10;
```

Search tables...

**Databases**

banala
banalas
default
foodmart
shyni
xademo

Execute    Explain    Save as...    New Worksheet

100%

**Query Process Results (Status: SUCCEEDED)**    Save results... ▾

Logs    Results

Filter columns...    previous    next

| name | votes |
|------|-------|
| The Shawshank Redemption | 2400000 |
| The Dark Knight | 2400000 |
| Inception | 2100000 |
| Pulp Fiction | 1900000 |

---

Hive    Query    Saved Queries    History    UDFs    Upload Table

**Visualization**

Data Visualization    Data Explorer    Maximum Row Count: 10000    OK

SQL

**COLUMNS**

Abc  name
Abc  votes
#    COUNT

**POSITIONAL**

▾ x    ▾ Abc  name    ✕
▾ y    ▾ Abc  votes    ✕
▾ col    drop a column here
▾ row    drop a column here

**MARKS**    🗑 Clear

▾ size    drop a column here
▾ color    drop a column here
▾ shape    drop a column here
▾ detail    drop a column here
▾ text    drop a column here

TEZ

➢ From the above visualization I can conclude that the total votes polled for each movie were revelaed .

➢ Here the top 3  votes polled for the  movies are 1.2400000 (The shawshank redemption)
2.2400000 (The dark knight)
3.2100000(Inception)

➢ Total votes polled means 2400000 people watched the 2 different   movies .

8.

```
1 select director,count(*) AS TOTALMOVIES from moviesdata group by director order by TOTALMOVIES desc;
```

**Query Process Results (Status: SUCCEEDED)**

| director | totalmovies |
| --- | --- |
| Woody Allen | 38 |
| Clint Eastwood | 31 |
| Directors | 28 |
| Steven Spielberg | 27 |
| Ron Howard | 24 |

9.

**10.**

11.

```
1 count (*) as TOTALMOVIES, country from moviesdata group by country order by TOTALMOVIES desc limit 5;
```

SQL

TEZ

Search tables...

**Databases**

| | |
|---|---|
| banala | |
| banalas | |
| default | |
| moviesdata | ☰ |
| name | STRING |
| rating | STRING |
| genre | STRING |
| year | INT |
| released | STRING |
| score | DOUBLE |
| votes | INT |
| director | STRING |
| writer | STRING |
| star | STRING |
| country | STRING |
| budget | INT |
| gross | INT |
| company | STRING |
| runtime | INT |
| sample_07 | ☰ |
| sample_08 | ☰ |

Execute   Explain   Save as...                                                New Worksheet

**Query Process Results (Status: SUCCEEDED)**                                Save results... ▼

Logs    Results

Filter columns...                                                    previous   next

| totalmovies | country |
|---|---|
| 5475 | United States |
| 816 | United Kingdom |
| 279 | France |
| 190 | Canada |
| 117 | Germany |

> From the above visualisation I can conclude that  each country how many movies were watched will be revealed .

> Here the top 3 countries total watched movies are   1. United states (5475)

2. United kingdom (816)

3.france (279)

> That means from united states 5475 movies generated the top most  revenue.

12.

default

Search tables...

Databases

banala
banalas
default
foodmart
shyni
xademo

Worksheet ✖    capstone 3 ✖    capstone 4 ✖    capstone 13 ✖    capstone 5 ✖    capstone 12 ✖

```
1  select name ,budget from moviesdata group by budget,name order by budget desc limit 10;
```

SQL

TEZ

6

Execute    Explain    Save as...        New Worksheet

**Query Process Results (Status: SUCCEEDED)**      Save results... ▾

Logs    Results

Filter columns...        previous   next

| name | budget |
|---|---|
| Avengers: Endgame | 356000000 |
| Avengers: Infinity War | 321000000 |
| Star Wars: Episode VIII - The Last Jedi | 317000000 |
| Justice League | 300000000 |

13.

capstone 5 ✖   capstone 9 ✖   capstone 10 ✖   capstone 11 ✖   capstone 12 ✖

```
1  select name ,gross from moviesdata group by gross,name order by gross desc limit 10;
```

SQL

TEZ

13

Execute   Explain   Save as...

New Worksheet
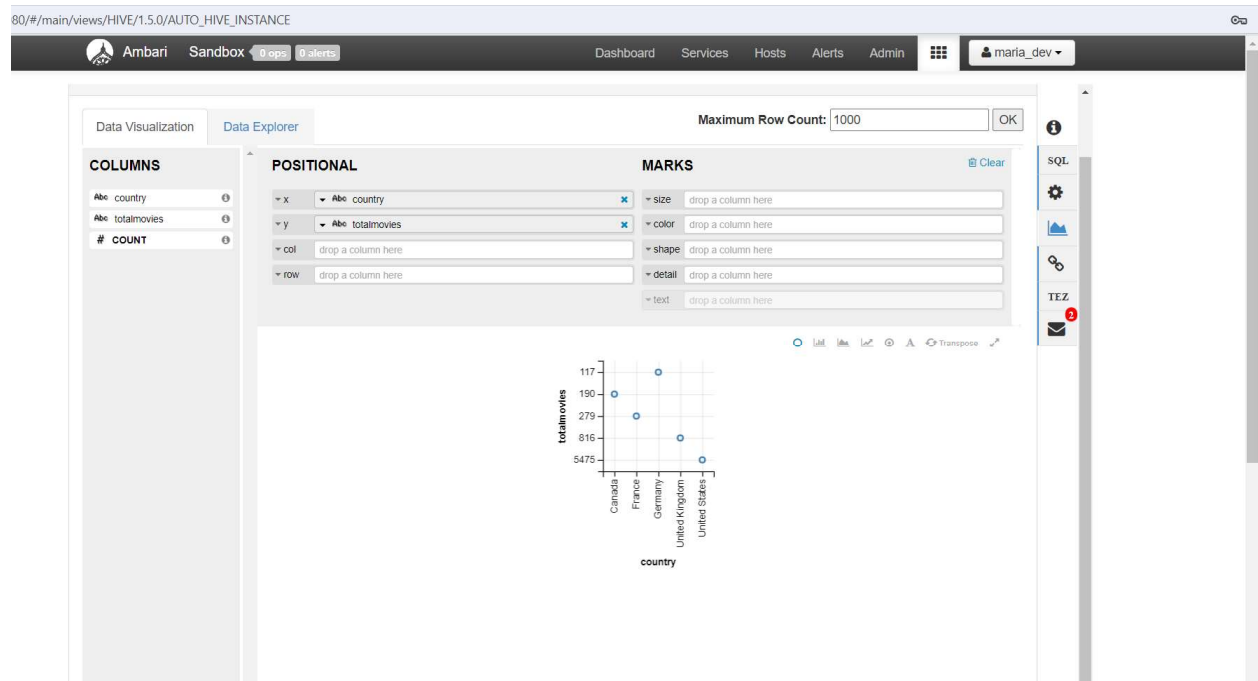
**Query Process Results (Status: SUCCEEDED)**

Save results... ▾

Logs   Results

Filter columns...

previous   next

| name | gross |
|---|---|
| Star Wars: Episode VII - The Force Awakens | 2069521700 |
| Avengers: Infinity War | 2048359754 |
| The Lion King | 1670727580 |
| Jurassic World | 1670516444 |
| The Avengers | 1518815515 |

14.

**Database Explorer**   ⟳

**Query Editor**   ⤢

default   ▾

capstone 5 ✖   capstone 9 ✖   capstone 10 ✖   capstone 11 ✖   capstone 13 ✖   Worksheet (14) ✖

```
1  select name,writer,gross from moviesdata order by gross desc limit 1;
```

Search tables...

Databases

banala
banalas
default
  moviesdata
  name        STRING
  rating      STRING
  genre       STRING
  year        INT
  released    STRING
  score       DOUBLE
  votes       INT
  director    STRING
  writer      STRING
  star        STRING
  country     STRING
  budget      INT
  gross       INT
  company     STRING
  runtime     INT
  sample 07
  sample 08

SQL

TEZ

15

Execute   Explain   Save as...

New Worksheet

**Query Process Results (Status: SUCCEEDED)**

Save results... ▾

Logs   Results

Filter columns...

previous   next

| name | writer | gross |
|---|---|---|
| Star Wars: Episode VII - The Force Awakens | Lawrence Kasdan | 2069521700 |

15.



16.

default

Search tables...

Databases

banala
banalas
default
  moviesdata
    name          STRING
    rating        STRING
    genre         STRING
    year          INT
    released      STRING
    score         DOUBLE
    votes         INT
    director      STRING
    writer        STRING
    star          STRING
    country       STRING
    budget        INT
    gross         INT
    company       STRING
    runtime       INT
  sample  07
  sample  08

capstone 5 ✖   capstone 9 ✖   capstone 10 ✖   capstone 11 ✖   capstone 13 ✖   capstone 14 ✖

capstone 15 ✖   Worksheet (17) ✖   Worksheet (18) * ✖

```
1  select name, gross/budget as RATIO from moviesdata order by RATIO desc limit 10;
```

Execute    Explain    Save as...                    New Worksheet

**Query Process Results (Status: SUCCEEDED)**                    Save results... ▾

Logs    Results

Filter columns...                                     previous    next

| name | ratio |
| --- | --- |
| Paranormal Activity | 12890.386666666667 |
| The Blair Witch Project | 4143.984983333334 |
| The Gallows | 429.6441 |

---

    votes         INT
    director      STRING
    writer        STRING
    star          STRING
    country       STRING
    budget        INT
    gross         INT
    company       STRING
    runtime       INT
  sample  07
  sample  08

Execute    Explain    Save as...                    New Worksheet

**Query Process Results (Status: SUCCEEDED)**                    Save results... ▾

Logs    Results

Filter columns...                                     previous    next

| name | ratio |
| --- | --- |
| Paranormal Activity | 12890.386666666667 |
| The Blair Witch Project | 4143.984983333334 |
| The Gallows | 429.6441 |
| El Mariachi | 291.56 |
| Once | 139.57814666666667 |
| Clerks | 116.70851851851852 |
| Napoleon Dynamite | 115.3472175 |
| In the Company of Men | 112.17892 |
| Keeping Mum | 109.98126627218934 |
| Open Water | 109.366974 |

➢ From the above visualization I can conclude that  gross to budget ratio means the profit  of each movie is revealed here

➢ That means top 3 highest profitable movies   are 1.Paranormal activity  (12890.38)
                                                     2.The blair witch project (4143.98)
                                                     3. The gallows (429.64)

17.

capstone 5 ✖   capstone 9 ✖   capstone 10 ✖   capstone 11 ✖   capstone 13 ✖   capstone 14 ✖

capstone 15 ✖   capstone 17 ✖   Worksheet (18) ✖

```
1  select count(*) as MAXMOVIES,company from moviesdata group by company order by MAXMOVIES desc limit 1
```

SQL

TEZ

**Databases**

banala
banalas
default
  moviesdata
    name        STRING
    rating      STRING
    genre       STRING
    year        INT
    released    STRING
    score       DOUBLE
    votes       INT
    director    STRING
    writer      STRING
    star        STRING
    country     STRING
    budget      INT
    gross       INT
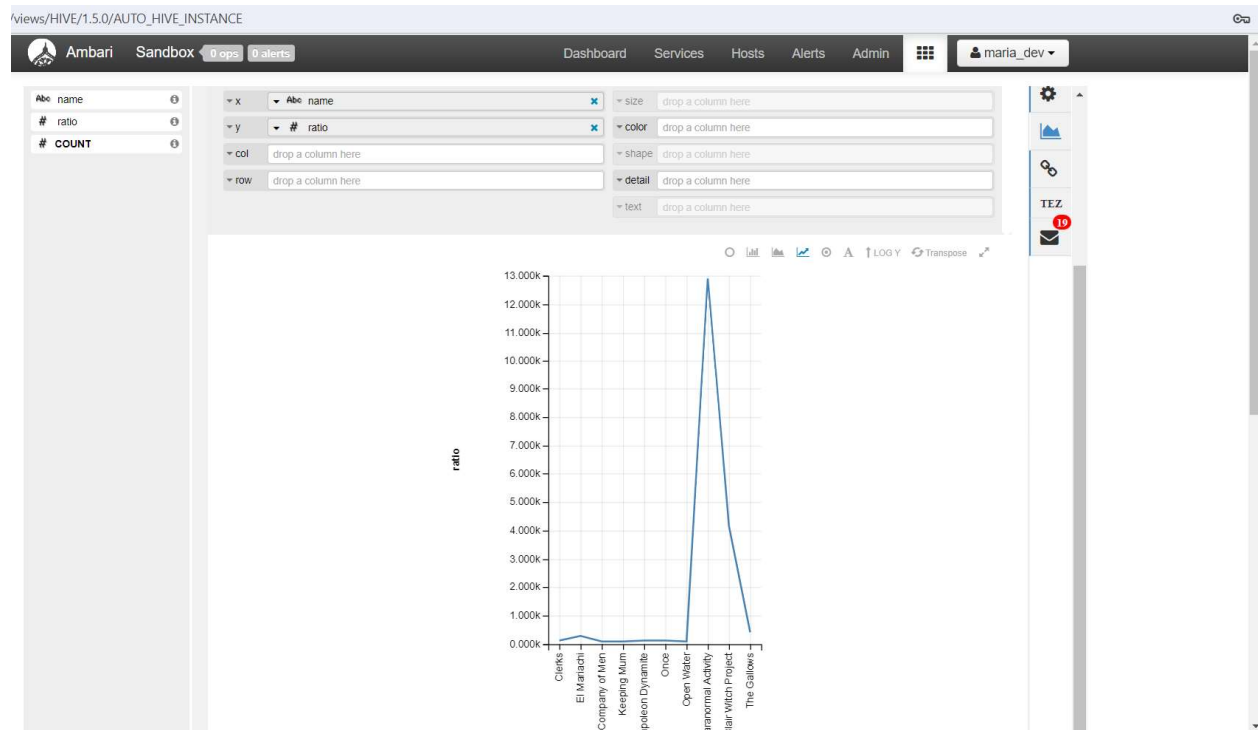    company     STRING
    runtime     INT
  sample 07
  sample 08

Search tables...

[Execute]  [Explain]  [Save as...]                    [New Worksheet]

100%

**Query Process Results (Status: SUCCEEDED)**        Save results... ▾

[Logs]  [Results]

Filter columns...                          [previous] [next]

| maxmovies | company |
|---|---|
| 377 | Universal Pictures |
| 334 | Warner Bros. |
| 332 | Columbia Pictures |

---

    year        INT
    released    STRING
    score       DOUBLE
    votes       INT
    director    STRING
    writer      STRING
    star        STRING
    country     STRING
    budget      INT
    gross       INT
    company     STRING
    runtime     INT
  sample 07
  sample 08

[Execute]  [Explain]  [Save as...]                    [New Worksheet]
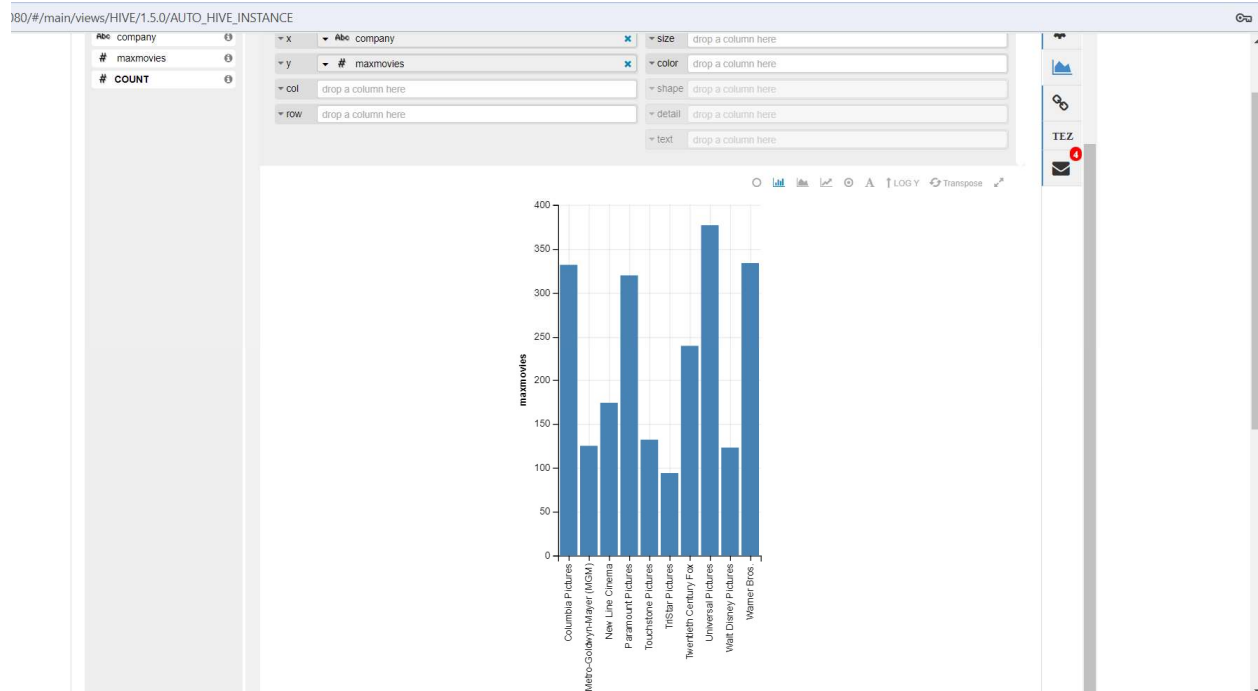
100%

**Query Process Results (Status: SUCCEEDED)**        Save results... ▾

[Logs]  [Results]

Filter columns...                          [previous] [next]

| maxmovies | company |
|---|---|
| 377 | Universal Pictures |
| 334 | Warner Bros. |
| 332 | Columbia Pictures |
| 320 | Paramount Pictures |
| 240 | Twentieth Century Fox |
| 174 | New Line Cinema |
| 132 | Touchstone Pictures |
| 125 | Metro-Goldwyn-Mayer (MGM) |
| 123 | Walt Disney Pictures |
| 94 | TriStar Pictures |

> ➤ From the above visualisation I can conclude that which company make the highest number of movies were revealed here.
> ➤ Here the top 3 companies make the movies are 1.universal pictures (377)
>                                              2.warner bros (334)
>                                              3.columbia pictures (332)

> ➤ From the above each company total movies production is revelaed .

18.

```
1  select name,runtime  from moviesdata order by runtime desc;
```

TEZ 7

**Execute**  Explain  Save as...

**New Worksheet**

**Query Process Results (Status: SUCCEEDED)**

Save results... ▼

Logs  **Results**

Filter columns...

previous  next

| name | runtime |
| --- | --- |
| The Best of Youth | 366 |
| Little Dorrit | 357 |
| Gettysburg | 271 |
| Hamlet | 242 |
| The Beautiful Troublemaker | 238 |

Databases

| rating | STRING |
| --- | --- |
| genre | STRING |
| year | INT |
| released | STRING |
| score | DOUBLE |
| votes | INT |
| director | STRING |
| writer | STRING |
| star | STRING |
| country | STRING |
| budget | INT |
| gross | INT |
| company | STRING |
| runtime | INT |

sample 07
sample 08
n
n2
n22
nu
foodmart

---

Search tables...

capstone 15 ✖  capstone 17 ✖  capstone 17 ✖  capstone 18 a ✖  Worksheet (20) ✖  Worksheet (21) ✖

SQL

```
1  select name,runtime  from moviesdata order by runtime;
```

TEZ 9

**Execute**  Explain  Save as...

**New Worksheet**

**Query Process Results (Status: SUCCEEDED)**

Save results... ▼

Databases

| rating | STRING |
| --- | --- |
| genre | STRING |
| year | INT |
| released | STRING |
| score | DOUBLE |
| votes | INT |
| director | STRING |
| writer | STRING |
| star | STRING |
| country | STRING |
| budget | INT |
| gross | INT |
| company | STRING |
| runtime | INT |

sample 07
sample 08
n
n2
n22
nu
foodmart

Logs  **Results**

Filter columns...

previous  next

| name | runtime |
| --- | --- |
| One for the Money | null |
| The Wolfman | null |
| Saving Mbango | null |
| Saw: The Final Chapter | null |
| The Business of Show Business | 55 |

19.

a.

capstone 5 ✖   capstone 10 ✖   capstone 17 ✖   capstone 18 a ✖   capstone 18 b ✖   capstone 19 a ✖

capstone 19 b ✖   capstone 19 c ✖   capstone 19 d * ✖

```
1  select avg(score) AS AVGSCORE from moviesdata;
```

Execute   Explain   Save as...                                   New Worksheet

100%

**Query Process Results (Status: SUCCEEDED)**                    Save results... ▼

Logs   Results

Filter columns...                                        previous   next

**avgscore**

6.390410958904098

---

b.

capstone 5 ✖   capstone 10 ✖   capstone 17 ✖   capstone 18 a ✖   capstone 18 b ✖   capstone 19 a ✖

capstone 19 b ✖   capstone 19 c ✖   capstone 19 d * ✖

```
1  select avg(budget) from moviesdata;
```

Execute   Explain   Save as...                                   New Worksheet

**Query Process Results (Status: SUCCEEDED)**                    Save results... ▼

Logs   Results

Filter columns...                                        previous   next

**_c0**

3.5589876192650534E7

c.

select avg(gross) AS AVGGROSS from moviesdata;

Query Process Results (Status: SUCCEEDED)

avggross

7.748249752233815E7

d.

default ▾

Search tables...

**Databases**

- banala
- banalas
- default
  - moviesdata ≡

| | |
|---|---|
| name | STRING |
| rating | STRING |
| genre | STRING |
| year | INT |
| released | STRING |
| score | DOUBLE |
| votes | INT |
| director | STRING |
| writer | STRING |
| star | STRING |
| country | STRING |
| budget | INT |
| gross | INT |
| company | STRING |
| runtime | INT |

- sample_07 ≡
- sample_08 ≡

```
1  select avg(runtime) AS AVGRUNTIME from moviesdata;
```

Execute  Explain  Save as...                    New Worksheet

**Query Process Results (Status: SUCCEEDED)**                 Save results... ▾

Logs  Results

Filter columns...                              previous  next

**avgruntime**

107.2616127348643

SQL
⚙
📊
🔗
TEZ ①
✉