

CSP:554 – Project Presentation

Group name: **Panda**

Arinjay Jain: **A20447307**

Ayush Dadhich: **A20449379**

Parth Gupta: **A20449774**

Keerthana Nagula: **A20442972**

- **Project Topic:** Develop a Data handling pipeline with the aid of Big Data SQL tools
- **Application subject area:** Movies & Entertainment
- **Data set Source:** IMDB Movies dataset from Kaggle
Data set of 1,000 most popular movies on IMDB (Tomiandrep, “IMDB Filmid”) in the last 10 years. The data points included are: Title, Genre, Description, Director, Actors, Year, Runtime, Rating, Votes, Revenue, Metascore.

The main Movies Metadata file (Banik, “The Movies Dataset”). Files contain information on 45,000 movies utilized in the full MovieLens dataset. Data set include plot keywords, languages, production countries and represents a collection of 26 million ratings from all the 45,000 movies.

- **Project Motive / Question:**
Develop a data processing pipeline using Big Data technologies to insert, modify, transform and implement Big Data techniques to research relation between proposed data and derive meaningful / relevant understanding of data streaming on movies data from IMDB.
- **Proposed Approach: Tools / Techniques for:**
 - **Data Insertion:** HDFS commands to move data on to Hadoop environment
 - **Data Modification & Transformation:** SparkSQL
 - **Database:** NoSQL HBASE DB
 - **Data Mining & Analysis:** SQL and Python / R on Big Data
- **Reference resources:**

Tomiandrep. “IMDB Filmid.” *Kaggle*, Kaggle, 13 Dec. 2017,

<https://www.kaggle.com/tomiandrep/imdb-filmid/data>.

Banik, Rounak. “The Movies Dataset.” *Kaggle*, 10 Nov. 2017,

<https://www.kaggle.com/rounakbanik/the-movies-dataset/data>.

Creating Hbase Cluster:-

portal.azure.com/#create/Microsoft.HDInsightCluster

AppsPython | KaggleTrending - Mo...New Tab

Microsoft AzureSearch resources, services, and docs (G+)

Home > New > Create HDInsight cluster

Create HDInsight cluster

[Go to classic create experience](#)

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *

Azure for Students

Resource group *

my_sandbox

[Create new](#)

Cluster details

Name your cluster, pick a region, and choose a cluster type and version. [Learn more](#)

Cluster name *

hbase

Region *

Central US

Cluster type *

HBase

[Change](#)

Version *

HBase 2.0.0 (HDI 4.0)

Cluster credentials

Enter new credentials that will be used to administer or access the cluster.

Cluster login username * ⓘ

admin

Cluster login password *

.....

Confirm cluster login password *

.....

Secure Shell (SSH) username * ⓘ

sshuser

Use cluster login password for SSH

☒

Review + create

« Previous

Next: Storage »

Creating Spark Cluster:- 1.

[←](#) [→](#) [🏠](#) [portal.azure.com/#create/Microsoft.HDInsightCluster](#)

Apps [Python | Kaggle](#) [🌐 Trending - Mo...](#) [🔍 New Tab](#)

Microsoft Azure [🔍 Search](#)

[Home](#) > [New](#) > [Create HDInsight cluster](#)

Create HDInsight cluster

[↗️ Go to classic create experience](#)

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *

Azure for Students

Resource group *

my_sandbox

[Create new](#)

Cluster details

Name your cluster, pick a region, and choose a cluster type and version. [Learn more](#)

Cluster name *

mysparkk

✓

Region *

Central US

▼

Cluster type *

Spark
[Change](#)

Version *

Spark 2.3 (HDI 3.6)

▼

Cluster credentials

Enter new credentials that will be used to administer or access the cluster.

Cluster login username * ⓘ

spark

✓

Cluster login password *

.....

✓

Confirm cluster login password *

.....

✓

Secure Shell (SSH) username * ⓘ

sshuser

✓

Use cluster login password for SSH

☒

[Review + create](#)

[« Previous](#)

[Next: Storage »](#)

2.

portal.azure.com/#create/Microsoft.HDInsightCluster

Apps Python | Kaggle Trending - Mo... New Tab

Microsoft Azure

Home > New > Create HDInsight cluster

Create HDInsight cluster

[Go to classic create experience](#)

Select or create a storage account that will be the default location for cluster logs and other output.

Primary storage type *

Azure Storage

Selection method * ⓘ

☒ Select from list ☐ Use access key

Primary storage account *

mysandboxdiag946

[Create new](#)

Container * ⓘ

mysparkk-2019-11-18t01-58-41-914z

✓

Data Lake Storage Gen1

Provide details for the cluster to access Data Lake Storage Gen1. The cluster will be able to access any Data Lake Storage Gen1 accounts that the chosen service principal has access to.

Data Lake Storage Gen1 access

[Configure access settings](#)

Additional Azure storage

Link additional Azure storage accounts to the cluster.

Account name

[Add Azure storage](#)

Metastore settings

To preserve your Hive and/or Oozie metadata outside of this cluster, select a SQL database for this cluster.

SQL database for Hive ⓘ

SQL database for Oozie ⓘ

Review + create

« Previous

Next: Security + networking »

3.

← → ↺ 🏠

portal.azure.com/#create/Microsoft.HDInsightCluster

🌐 Apps k Python | Kaggle 🌐 Trending - Mo... 🌐 New Tab

Microsoft Azure 🔍 Search

Home > New > Create HDInsight cluster

Create HDInsight cluster

↔️ Go to classic create experience

✓ Validation succeeded.

networking, or data transfer.

Basics

Subscription	Azure for Students
Resource group	my_sandbox
Region	Central US
Cluster name	(new) mysparkk
Cluster type	Spark 2.3 (HDI 3.6)
Cluster login username	spark
Secure Shell (SSH) username	sshuser
Use cluster login password for SSH	Enabled

Security + networking

Virtual network	my_sandbox-vnet
Subnet	default

Storage

Primary storage type	Azure Storage
Primary storage account	mysandboxdiag946
Container	mysparkk-2019-11-18t01-58-41-914z
Additional Azure storage	None
Data Lake Storage Gen1 access	Disabled

Cluster configuration

Head	2 nodes, D12 v2 (4 Cores, 28 GB RAM)
Worker	4 nodes, D13 v2 (8 Cores, 56 GB RAM)

Create« PreviousNextDownload a template for automation

⏪

Create Table

```
bigdata@bigdata:~ x sshuser@hn0-hbase: ~ x ▼
Took 1.4901 seconds
hbase(main):003:0> drop 'imdb_mdata'
Took 4.3238 seconds
hbase(main):004:0> disable 'movie_mdata'
Took 1.2663 seconds
hbase(main):005:0> drop 'movie_mdata'
Took 2.2654 seconds
hbase(main):006:0> create 'imdb_mdata', 'cf'
Created table imdb_mdata
Took 4.2858 seconds
=> Hbase::Table - imdb_mdata
hbase(main):007:0> create 'movie_mdata', 'cf'
Created table movie_mdata
Took 4.3300 seconds
=> Hbase::Table - movie_mdata
hbase(main):008:0> list
TABLE
imdb_mdata
movie_mdata
2 row(s)
Took 0.0313 seconds
=> ["imdb_mdata", "movie_mdata"]
hbase(main):009:0> 
```

Load Data

```
sshuser@hn0-hbase: ~  
bigdata@bigdata:~ x sshuser@hn0-hbase: ~ x  
sshuser@hn0-hbase:~$ hbase org.apache.hadoop.hbase.mapreduce.ImportT  
sv -Dimporttsv.separator="," -Dimporttsv.columns="HBASE_ROW_KEY,cf:a  
dult,cf:budget,cf:genres,cf:imdb_id,cf:original_language,cf:original  
_title,cf:overview,cf:popularity,cf:production_companies,cf:producti  
on_countries,cf:release_date,cf:revenue,cf:runtime,cf:spoken_languag  
es,cf:status,cf:tagline,cf:vote_average,cf:vote_count" movie_mdata /  
movies_mdata.csv  
SLF4J: Class path contains multiple SLF4J bindings.  
SLF4J: Found binding in [jar:file:/usr/hdp/3.1.2.2-1/phoenix/phoenix  
-5.0.0.3.1.2.2-1-server.jar!/org/slf4j/impl/StaticLoggerBinder.class  
]  
SLF4J: Found binding in [jar:file:/usr/hdp/3.1.2.2-1/hadoop/lib/slf4  
j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]  
SLF4J: See http://www.slf4j.org/codes.html#multiple\_bindings for an  
explanation.  
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]  
2019-11-24 21:40:14,116 INFO [main] zookeeper.ReadOnlyZKClient: Con  
nect 0x6f45df59 to zk0-hbase.kk4rdbp1quurlqz3nexzcjo5g.gx.internal.  
cloudapp.net:2181,zk6-hbase.kk4rdbp1quurlqz3nexzcjo5g.gx.internal.c  
loudapp.net:2181,zk4-hbase.kk4rdbp1quurlqz3nexzcjo5g.gx.internal.cl  
oudapp.net:2181 with session timeout=120000ms, retries 6, retry inte  
rval 1000ms, keepAlive=60000ms
```


version 2.3.2.2.6.5.3009-43

```
scala> val df = withCatalog(catalog)
19/11/18 17:34:50 WARN Configuration: hbase-site.xml:an attempt to override final parameter: dfs.support.append; Ignoring.
df: org.apache.spark.sql.DataFrame = [rowkey: string, actors: string ... 10 more fields]

scala> df.show()
```

```
Q sshuser@hn0-myspar: ~
sshuser@hn0-hbase: ~
x sshuser@hn0-myspar: ~
19/11/18 17:36:26 WARN Configuration: hbase-site.xml:an attempt to override final parameter: dfs.support.append; Ignoring.
19/11/18 17:36:27 WARN Configuration: hbase-site.xml:an attempt to override final parameter: dfs.support.append; Ignoring.
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|rowkey|actors|description|director|genre|metascore|rating|revenue|runtime_min|title|votes|year|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|1|Chris Pratt; Vin ...|A group of interg...|James Gunn|Action;Adventure;...|76|8.1|333.13|121|Guardians of the ...|121|757074|2014|
|10|Jennifer Lawrence...|A spacecraft trav...|Morten Tyldum|Adventure;Drama;R...|41|7|100.01|116|Passengers|192177|2016|
|100|Leonardo DiCaprio...|An undercover cop...|Martin Scorsese|Crime;Drama;Thriller|85|8.5|132.37|151|The Departed|937414|2006|
|1000|Kevin Spacey; Jen...|A stuffy business...|Barry Sonnenfeld|Comedy;Family;Fan...|11|5.3|19.64|87|Nine Lives|12435|2016|
|101|Tom Hardy; Emily ...|Identical twin ga...|Brian Helgeland|Biography;Crime;D...|55|7|1.87|132|Legend|108836|2015|
|102|Chris Hemsworth; ...|The powerful but ...|Kenneth Branagh|Action;Adventure;...|57|7|181.02|115|Thor|570814|2011|
|103|Matt Damon; Jessi...|An astronaut beco...|Ridley Scott|Adventure;Drama;S...|80|8|228.43|144|The Martian|556097|2015|
|104|Marlo Casas; Ana ...|A young businessm...|Orlto Paulo|Crime;Mystery;Thr...|null|7.9|null|106|Contratempo|7204|2016|
|105|Henry Cavill; Arm...|In the early 1960...|Guy Ritchie|Action;Adventure;...|56|7.3|45.43|116|The Man from U.N....|202973|2015|
|106|Chris Pine; Ben F...|A divorced father...|David Mackenzie|Crime;Drama;Thriller|40|5.4|1.66|102|Hell or High Water|115546|2016|
|107|Robert De Niro; L...|A look at the lif...|Taylor Hackford|Comedy|40|5.4|1.66|120|The Comedian|1954|2016|
|108|Alexander Skarsgå...|Tarzan having acc...|David Yates|Action;Adventure;...|44|6.3|126.59|110|The Legend of Tarzan|117590|2016|
|109|Eve Lindley; Rich...|A mother struggle...|Katie Holmes|Drama|48|5.8|null|105|All We Had|1084|2016|
|11|Eddie Redmayne; K...|The adventures of...|David Yates|Adventure;Family;...|66|7.5|234.02|133|Fantastic Beasts ...|232072|2016|
|110|Alicia Vikander; ...|A young programme...|Alex Garland|Drama;Mystery;Sci-Fi|78|7.7|25.44|108|Ex Machina|339797|2014|
|111|John Gallagher Jr...|In a twisted soci...|Greg McLean|Action;Horror;Thr...|44|6.3|10.16|89|The Belko Experiment|3712|2016|
|112|Chiwetel Ejiofor; ...|In the antebellum...|Steve McQueen|Biography;Drama;H...|96|8.1|56.67|134|12 Years a Slave|486338|2013|
|113|Keanu Reeves; Jas...|A dystopian love ...|Ana Lily Amirpour|Romance;Sci-Fi|65|6.1|null|118|The Bad Batch|512|2016|
|114|Gerard Butler; Le...|King Leonidas of ...|Zack Snyder|Action;Fantasy;War|52|7.7|210.59|117|300|637104|2006|
|115|Dante Redcliffe; ...|Harry Ron and Her...|David Yates|Adventure;Drama;F...|87|8.1|380.96|130|Harry Potter and ...|590595|2011|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows

scala> df.createTempView("imdb_Mdata")
scala> spark.sqlContext.sql("select title, rating from imdb_Mdata").show
+-----+-----+
|title|rating|
+-----+-----+
|Guardians of the ...|8.1|
|Passengers|7|
|The Departed|8.5|
|Nine Lives|5.3|
|Legend|7|
|Thor|7|
|The Martian|8|
|Contratempo|7.9|
|The Man from U.N....|7.3|
|Hell or High Water|7.7|
|The Comedian|5.4|
|The Legend of Tarzan|6.3|
|All We Had|5.8|
|Fantastic Beasts ...|7.5|
|Ex Machina|7.7|
|The Belko Experiment|6.3|
|12 Years a Slave|8.1|
|The Bad Batch|6.1|
|300|7.7|
+-----+-----+
only showing top 20 rows
```

```
scala> def withCatalog(cat: String): DataFrame = {
  | spark.sqlContext
  | .read
  | .options(Map(HBaseTableCatalog.tableCatalog->cat))
  | .format("org.apache.spark.sql.execution.datasources.hbase")
  | .load()
  | }
withCatalog: (cat: String)org.apache.spark.sql.DataFrame

scala> val df = withCatalog(catalog)
19/11/18 21:51:01 WARN Configuration: hbase-site.xml:an attempt to override final parameter: dfs.support.append; Ignoring.
df: org.apache.spark.sql.DataFrame = [rowkey: string, adult: string ... 17 more fields]

scala> df.show()
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|id|adult|budget|genres|imdb_id|original_language|original_title|overview|popularity|production_companies|production_countries|release_date|revenue|runtime|spo
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|ken_languages|status|tagline|vote_average|vote_count|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|862|FALSE|30000000|[[{"id": 16 "name"...|t0114709|5415|en|Toy Story|Led by Woody Andy...|21.946943|[[{"name": 'Pixar ...|['iso_3166_1': '...' 10/30/1995|373554033|81|['iso_
639_1': 'e...|Released|Roll the dice and...|2413|en|Jumanji|When siblings Jud...|17.015539|[[{"name": 'TriSta...|['iso_3166_1': '...' 12/15/1995|262797249|104|['iso_
8844|FALSE|65000000|[[{"id": 12 "name"...|t0113497|6.9|en|Grumpier Old Men|A family wedding ...|11.7129|[[{"name": 'Warner...|['iso_3166_1': '...' 12/22/1995|0|101|['iso_
639_1': 'e...|Released|Still feeling St...|6.5|92|en|Waiting to Exhale|Cheated on mistr...|3.859495|[[{"name": 'Twenti...|['iso_3166_1': '...' 12/22/1995|81452156|127|['iso_
13157|FALSE|16000000|[[{"id": 35 "name"...|t0114885|6.1|34|en|Father of the Bri...|Just when George ...|8.387519|[[{"name": 'Sando...|['iso_3166_1': '...' 2/10/1995|76578911|106|['iso_
11862|FALSE|0|[[{"id": 35 "name"...|t0113041|5.7|173|en|Heat|Obsessive master ...|17.924927|[[{"name": 'Regenc...|['iso_3166_1': '...' 12/15/1995|187436818|170|['iso_
639_1': 'e...|Released|Just When His Wor...|5.7|1886|en|Sabrina|An ugly duckling ...|6.677277|[[{"name": 'Paramo...|['iso_3166_1': '...' 12/15/1995|0|127|['iso_
949|FALSE|60000000|[[{"id": 28 "name"...|t0113277|6.2|141|en|Tom and Huck|A mischievous you...|2.561161|[[{"name": 'Walt D...|['iso_3166_1': '...' 12/22/1995|0|97|['iso_
639_1': 'e...|Released|The original Bad ...|5.4|45|en|Sudden Death|International act...|5.23158|[[{"name": 'Univer...|['iso_3166_1': '...' 12/22/1995|64350171|106|['iso_
9991|FALSE|35000000|[[{"id": 28 "name"...|t0114576|5.5|174|en|GoldenEye|James Bond must u...|14.686036|[[{"name": 'Unlited...|['iso_3166_1': '...' 11/16/1995|352194034|130|['iso_
639_1': 'e...|Released|No limits. No fea...|6.6|1194|en|The American Pres...|Widowed U.S. pres...|6.318445|[[{"name": 'Columb...|['iso_3166_1': '...' 11/17/1995|107879496|106|['iso_
9007|FALSE|62000000|[[{"id": 35 "name"...|t0112346|199|en|Dracula: Dead and...|When a lawyer sho...|5.430331|[[{"name": 'Columb...|['iso_3166_1': '...' 12/22/1995|0|88|['iso_
12110|FALSE|0|[[{"id": 35 "name"...|t0112896|5.7|210|en|Balto|An outcast half-w...|12.140733|[[{"name": 'Unliver...|['iso_3166_1': '...' 12/22/1995|11348324|78|['iso_
21032|FALSE|0|[[{"id": 10751 "na...|t0112453|423|en|Nixon|An all-star cast ...|5.092|[[{"name": 'Hollyw...|['iso_3166_1': '...' 12/22/1995|13681765|192|['iso_
10858|FALSE|44000000|[[{"id": 36 "name"...|t0113907|7.1|72|en|Cutthroat Island|Morgan Adams and ...|7.284447|[[{"name": 'Le Stu...|['iso_3166_1': '...' 12/22/1995|10017322|119|['iso_
1408|FALSE|98000000|[[{"id": 28 "name"...|t0112760|5.7|137|en|Castno|The life of the g...|10.137389|[[{"name": 'Univer...|['iso_3166_1': '...' 11/22/1995|116112375|178|['iso_
639_1': 'e...|Released|Lose your heart a...|7.2|364|en|Sense and Sensib...|Rich Mr. Dashwood...|10.673167|[[{"name": 'Columb...|['iso_3166_1': '...' 12/13/1995|135000000|136|['iso_
514|FALSE|52000000|[[{"id": 18 "name"...|t0112641|7.8|539|en|Four Rooms|It's Ted the Bell...|9.026586|[[{"name": 'Mirana...|['iso_3166_1': '...' 12/9/1995|4300000|90|['iso_
639_1': 'e...|Released|Twelve outrageous...|6.5|1128|en|Ace Ventura: When...|Summoned from an ...|8.205448|[[{"name": 'O Ente...|['iso_3166_1': '...' 11/10/1995|212385533|98|['iso_
9273|FALSE|30000000|[[{"id": 80 "name"...|t0112281|6.1|
639_1': 'e...|Released|New animals. New ...|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
>>> |
```

```
-----movie count each release year-----
```

-----Movie count each director directed-----

-----movie count each year with rating-----

-----Revenue Stylized Facts per year-----

-----Avg. Revenue of movies per director-----

```

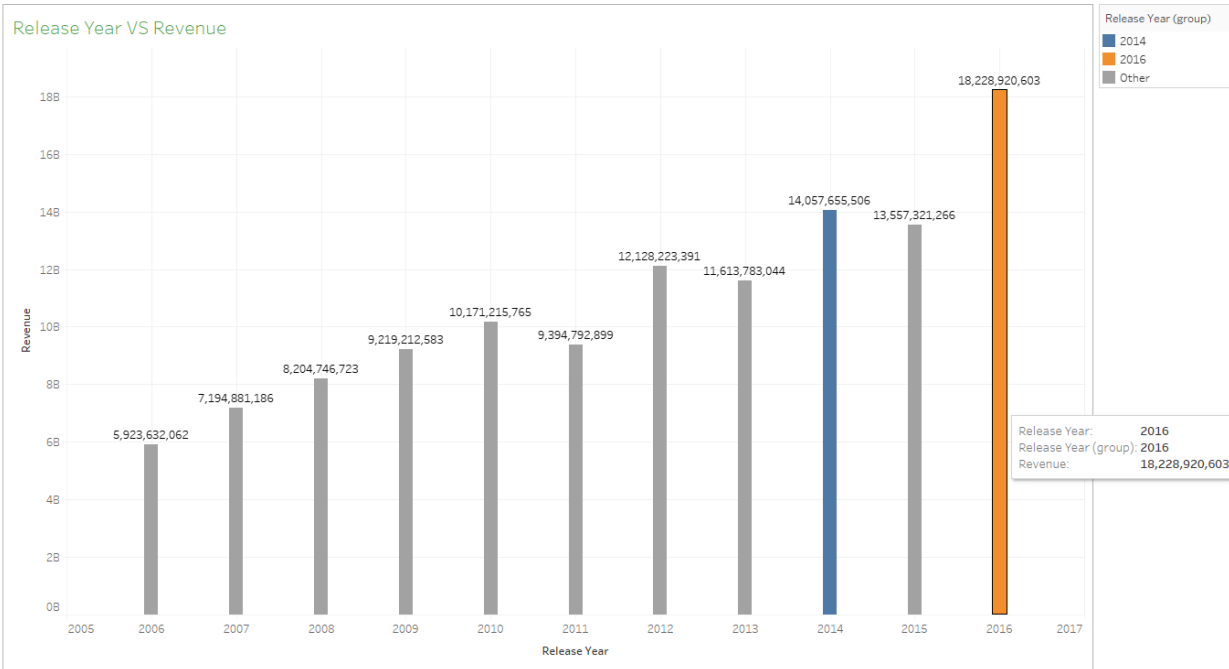
>>> result8 = spark.sql("SELECT c4, avg_c10) as Avg_revenue from df group by c4 HAVING Avg_revenue > 100 order by Avg_revenue DESC")
>>> result8.collect()
[Stage 33:=====>] (126 + 4) [Stage 33:=====>] (177 + 5)
[Stage 35:=====>] (157 + 4) [Stage 35:=====>] (80 + 5) [Stage 35:=====>] (118 + 4) [Stage 35:=====>]
=====
(157 + 4) [Stage 35:=====>] (80 + 5) [Stage 35:=====>] (118 + 4) [Stage 35:=====>]
=====
[Row(c4='Joss Whedon', Avg_revenue=541.13499999999999), Row(c4='Lee Unkrich', Avg_revenue=414.98000000000002), Row(c4='Gary Ross', Avg_revenue=488.0), Row(c4='Chris Buck', Avg_revenue=480.74000000000001), Row(c4='Chris Renaud', Avg_revenue=308.31), Row(c4='Gareth Edwards', Avg_revenue=366.41999999999996), Row(c4='Tim Miller', Avg_revenue=363.01999999999998), Row(c4='Byron Howard', Avg_revenue=341.25999999999999), Row(c4='J.J. Abrams', Avg_revenue=336.68999999999994), Row(c4='Kyle Balda', Avg_revenue=336.02999999999997), Row(c4='Anthony Russo', Avg_revenue=33.914999999999996), Row(c4='Francis Lawrence', Avg_revenue=324.95249999999999), Row(c4='Pete Docter', Avg_revenue=324.71500000000003), Row(c4='Pierre Coffin', Avg_revenue=309.77499999999998), Row(c4='Christopher Nolan', Avg_revenue=303.01800000000003), Row(c4='David Slade', Avg_revenue=306.51999999999998), Row(c4='Bill Condon', Avg_revenue=286.78999999999996), Row(c4='Sam Raimi', Avg_revenue=284.14999999999997), Row(c4='John Yates', Avg_revenue=271.6666666666667), Row(c4='Stephen Herek', Avg_revenue=270.99999999999999), Row(c4='Dan Scanlan', Avg_revenue=269.00000000000001), Row(c4='Andrew Stanton', Avg_revenue=261.05333333333333), Row(c4='Favre', Avg_revenue=259.99999999999998), Row(c4='Robert Stronberg', Avg_revenue=254), Row(c4='Mark Andrews', Avg_revenue=248), Row(c4='Michael Bay', Avg_revenue=236.88666666666666), Row(c4='Shane Black', Avg_revenue=222.02), Row(c4='Don Hall', Avg_revenue=222.49000000000001), Row(c4='John Lasseter', Avg_revenue=217.75), Row(c4='Mark Osborne', Avg_revenue=215.40000000000001), Row(c4='Peter Jackson', Avg_revenue=215.11250000000001), Row(c4='Gore Verbinski', Avg_revenue=207.45499999999998), Row(c4='Nathan Greno', Avg_revenue=200.81), Row(c4='Dean DeBlois', Avg_revenue=197.19499999999999), Row(c4='Bryan Singer', Avg_revenue=196.43666666666664), Row(c4='Phil Lord', Avg_revenue=195.94333333333333), Row(c4='Zack Snyder', Avg_revenue=195.148), Row(c4='Catherine Hardwicke', Avg_revenue=191.44999999999999), Row(c4='Rich Moore', Avg_revenue=189.41), Row(c4='Marc Forster', Avg_revenue=185.36000000000001), Row(c4='Rob Marshall', Avg_revenue=184.53), Row(c4='Tim Johnson', Avg_revenue=177.34), Row(c4='Joe Johnston', Avg_revenue=176.63999999999999), Row(c4='Ron Clements', Avg_revenue=176.56), Row(c4='George Miller', Avg_revenue=175.81), Row(c4='Sam Mendes', Avg_revenue=175.77000000000001), Row(c4='Brad Bird', Avg_revenue=169.74000000000001), Row(c4='Genndy Tartakovsky', Avg_revenue=169.69), Row(c4='Todd McGehee', Avg_revenue=169.27000000000001), Row(c4='Sanjay Dahiya', Avg_revenue=160.15000000000001), Row(c4='John Lasseter', Avg_revenue=159.99999999999999), Row(c4='Gregory La Roca', Avg_revenue=166.66666666666667), Row(c4='Justin Lita', Avg_revenue=159.8), Row(c4='Chris East', Avg_revenue=174.75), Row(c4='James Cameron', Avg_revenue=167.00000000000001), Row(c4='Todd Phillips', Avg_revenue=160.16499999999999), Row(c4='Steven Spielberg', Avg_revenue=156.7475), Row(c4='Brad Peyton', Avg_revenue=155.18000000000001), Row(c4='Rupert Sanders', Avg_revenue=155.11000000000001), Row(c4='Alfonso Cuarón', Avg_revenue=154.685), Row(c4='Walt Dohrn', Avg_revenue=153.69), Row(c4='Brett Ratner', Avg_revenue=153.50999999999999), Row(c4='David Ayer', Avg_revenue=150.56999999999999), Row(c4='Stephen Sommers', Avg_revenue=150.16999999999999), Row(c4='Peyton Reed', Avg_revenue=149.435), Row(c4='Tom McGrath', Avg_revenue=148.34), Row(c4='Alan Taylor', Avg_revenue=148.04500000000002), Row(c4='Kenneth Branagh', Avg_revenue=144.24000000000001), Row(c4='Phyllida Lloyd', Avg_revenue=143.69999999999999), Row(c4='Carlos Saldanha', Avg_revenue=143.62), Row(c4='Alessandro Carloni', Avg_revenue=143.52000000000001), Row(c4='Paul Feig', Avg_revenue=141.95500000000001), Row(c4='Martin Campbell', Avg_revenue=141.80000000000001), Row(c4='Scott Derrickson', Avg_revenue=140.32999999999998), Row(c4='Rawnson Marshall Thurber', Avg_revenue=138.875), Row(c4='Tim Burton', Avg_revenue=138.505), Row(c4='Jonathan Liebesman', Avg_revenue=137.255), Row(c4='James Mangold', Avg_revenue=136.00000000000001), Row(c4='John Schofield', Avg_revenue=131.5), Row(c4='Alex Proyas', Avg_revenue=130.99999999999999), Row(c4='John Dahl', Avg_revenue=129.78999999999999), Row(c4='Michael McCurrie', Avg_revenue=126.99999999999999), Row(c4='Adam Rodriguez', Avg_revenue=126.67999999999999), Row(c4='Michael Patrick King', Avg_revenue=123.94999999999999), Row(c4='Tate Taylor', Avg_revenue=122.51000000000001), Row(c4='Seth Gordon', Avg_revenue=117.02), Row(c4='John Lee Hancock', Avg_revenue=117.34666666666668), Row(c4='F. Gary Gray', Avg_revenue=117.185), Row(c4='Gabriele Muccino', Avg_revenue=116.27000000000001), Row(c4='Mark Steven Johnson', Avg_revenue=115.8), Row(c4='Angela Jolie', Avg_revenue=115.99999999999999), Row(c4='Anne Fletcher', Avg_revenue=114.60999999999999), Row(c4='Seth MacFarlane', Avg_revenue=114.17), Row(c4='Roland Emmerich', Avg_revenue=114.11666666666667), Row(c4='James Gunn', Avg_revenue=113.73999999999999), Row(c4='Jason Reitman', Avg_revenue=113.65000000000001), Row(c4='Tony Gilroy', Avg_revenue=113.17), Row(c4='Jon Lucas', Avg_revenue=113.08), Row(c4='Quentin Tarantino', Avg_revenue=112.48), Row(c4='Guy Ritchie', Avg_revenue=111.74250000000001), Row(c4='Ryan Coogler', Avg_revenue=109.70999999999999), Row(c4='Adam McKay', Avg_revenue=109.535), Row(c4='Peter Billingsley', Avg_revenue=109.18000000000001), Row(c4='Judd Apatow', Avg_revenue=108.75333333333333), Row(c4='David O. Russell', Avg_revenue=108.05500000000001), Row(c4='Clay Kaytis', Avg_revenue=107.51000000000001), Row(c4='Noam Murro', Avg_revenue=106.5), Row(c4='Chris Columbus', Avg_revenue=106.33333333333333), Row(c4='Richard Linklater', Avg_revenue=105.54400000000001), Row(c4='Rupert Wyatt', Avg_revenue=105.185), Row(c4='Harald Zwart', Avg_revenue=103), Row(c4='Sean Anders', Avg_revenue=102.36499999999999), Row(c4='Peter Berg', Avg_revenue=102.26599999999999), Row(c4='Steven Soderbergh', Avg_revenue=102.16333333333333), Row(c4='Matt Reeves', Avg_revenue=100.23333333333333), Row(c4='Tom Hooper', Avg_revenue=100.09666666666668)]

```

-----Ratings Stylized Facts per year-----

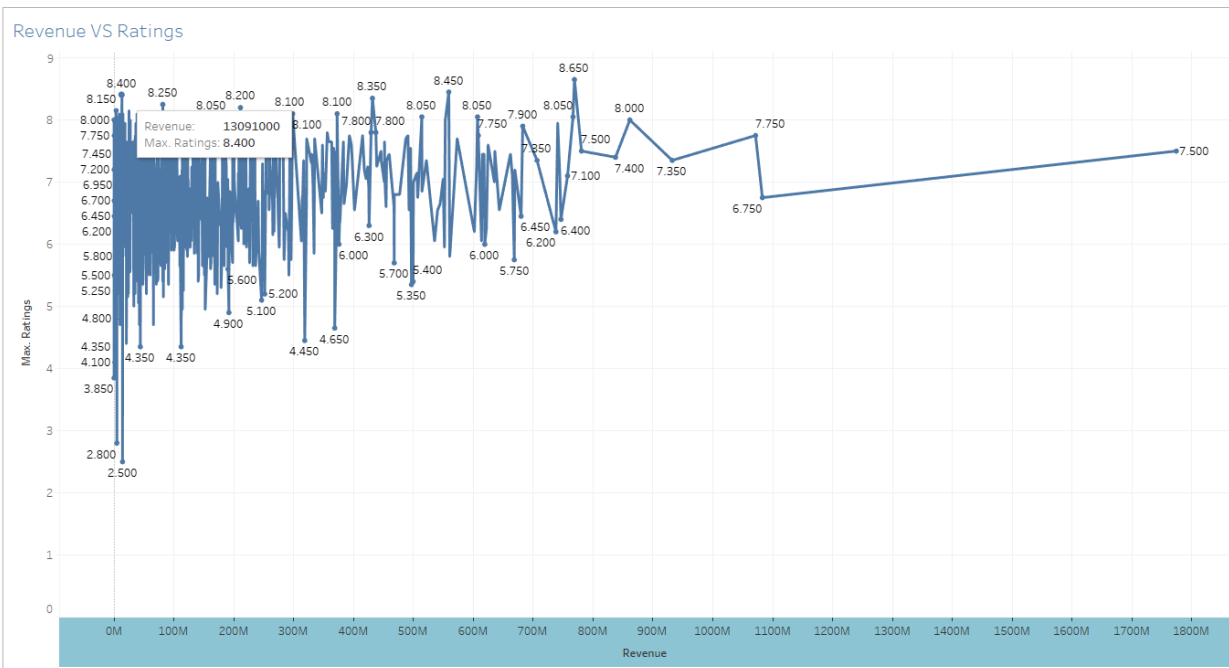
```
>>> result9 = spark.sql("SELECT _c6, sum(_c8) as sum_rating, avg(_c8) as avg_rating, min(_c8) as min_ratings, max(_c8) as max_rating from df group by _c6 order by _c6 DESC")
>>> result9.collect()
[Stage 38:=====] (160 + 4[Stage 38:=====] (197 + 3
[Stage 40:=====] (100 + 4[Stage 40:=====] (146 + 4[Stage 40:=====]
===== (191 + 4
Row(_c6=u'2015', sum_rating=838.50000000000023, avg_rating=6.602362204724411, min_ratings=u'3.5', max_rating=u'8.3'), Row(_c6=u'2014', sum_rating=670.10000000000002, avg_rating=6.8377551020408163, min_ratings=u'5.1', max_rating=u'8.6'), Row(_c6=u'2013', sum_rating=619.899999999999986, avg_rating=6.812087912087911, min_ratings=u'4.3', max_rating=u'8.2'), Row(_c6=u'2012', sum_rating=443.199999999999993, avg_rating=6.92499999999999989, min_ratings=u'5.3', max_rating=u'8.5'), Row(_c6=u'2011', sum_rating=430.800000000000007, avg_rating=6.8380952380952396, min_ratings=u'4.9', max_rating=u'8.6'), Row(_c6=u'2010', sum_rating=409.600000000000008, avg_rating=6.8266666666666668, min_ratings=u'4.2', max_rating=u'8.8'), Row(_c6=u'2009', sum_rating=355.0, avg_rating=6.9607843137254903, min_ratings=u'2.7', max_rating=u'8.4'), Row(_c6=u'2008', sum_rating=352.79999999999995, avg_rating=6.7846153846153836, min_ratings=u'1.9', max_rating=u'9'), Row(_c6=u'2007', sum_rating=378.099999999999991, avg_rating=7.1339622641509415, min_ratings=u'4.7', max_rating=u'8.5'), Row(_c6=u'2006', sum_rating=313.500000000000006, avg_rating=7.1250000000000009, min_ratings=u'5.6', max_rating=u'8.5'))]
>>>
```

Release Year VS Revenue



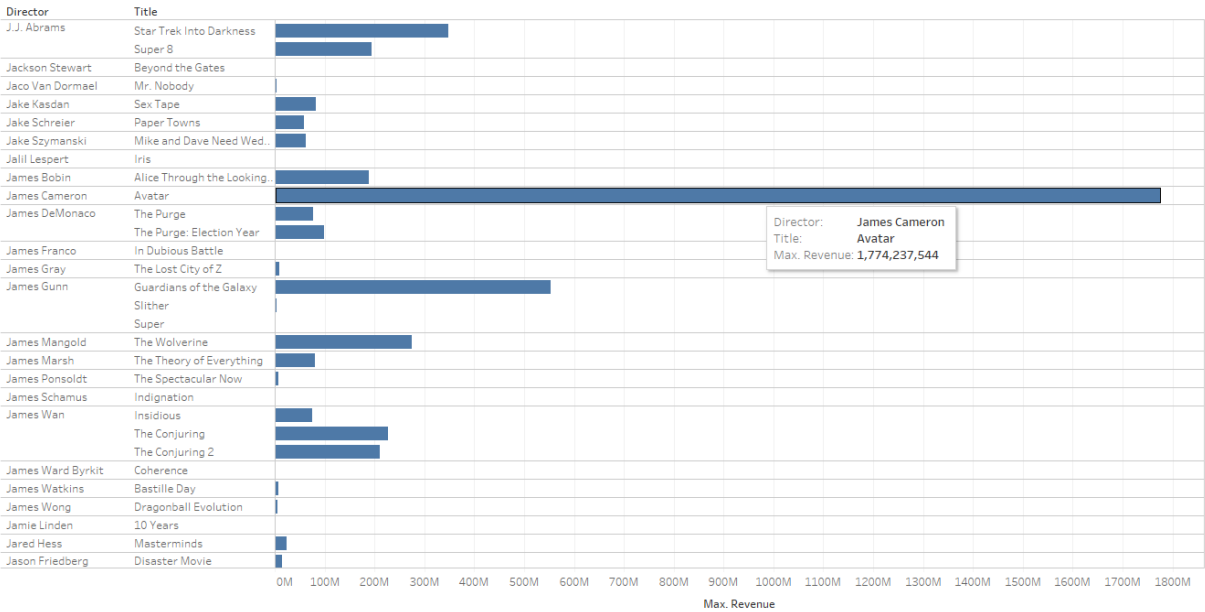
Sheet1 Sheet 2

Revenue VS Ratings



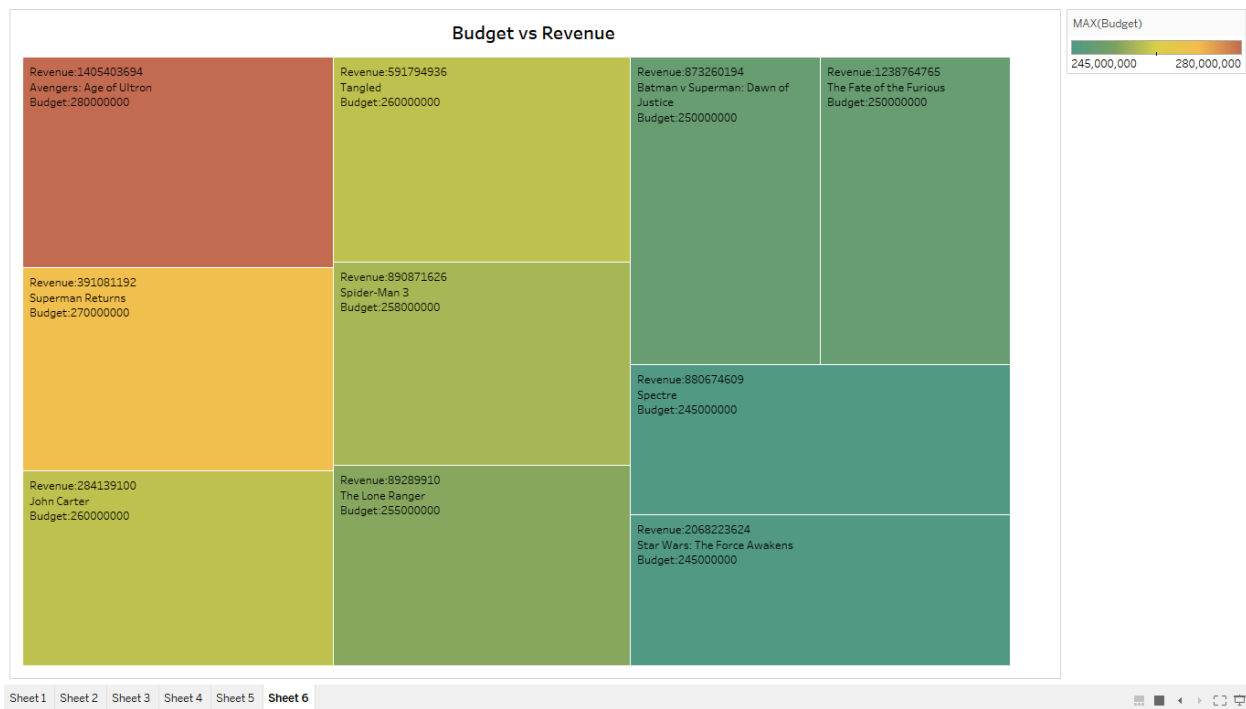
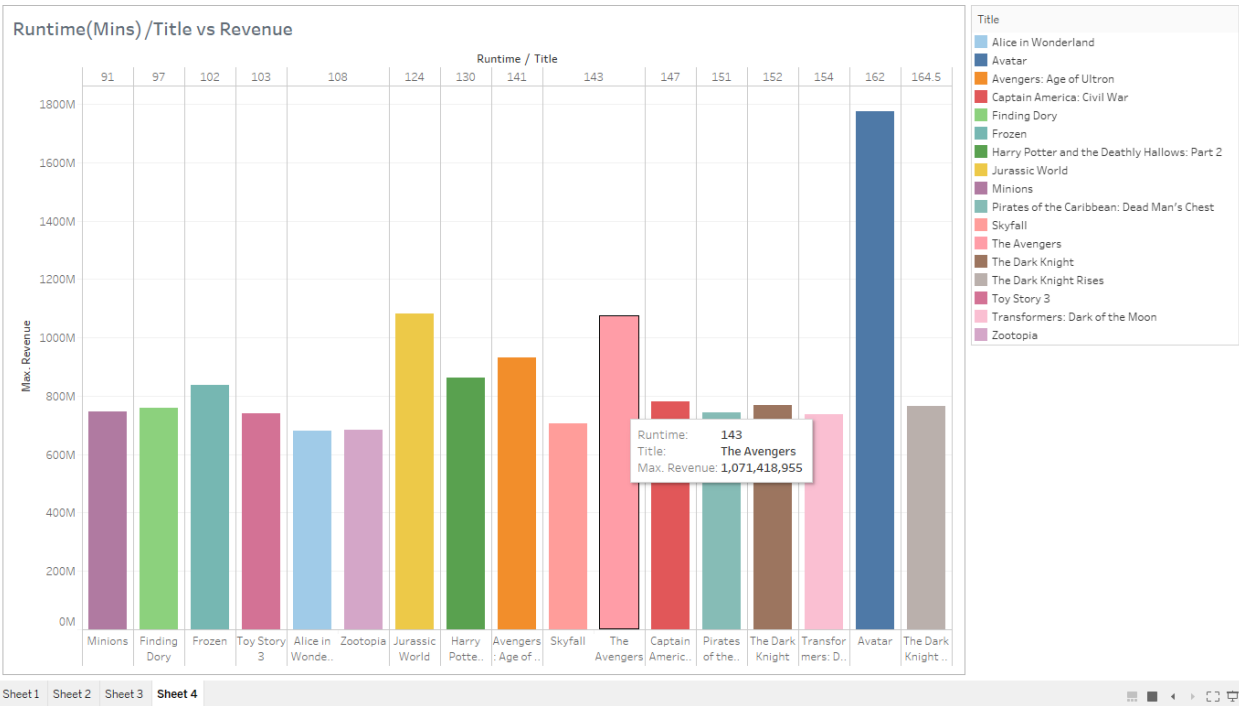
Sheet1 Sheet 2

Director, Movie and Revenue

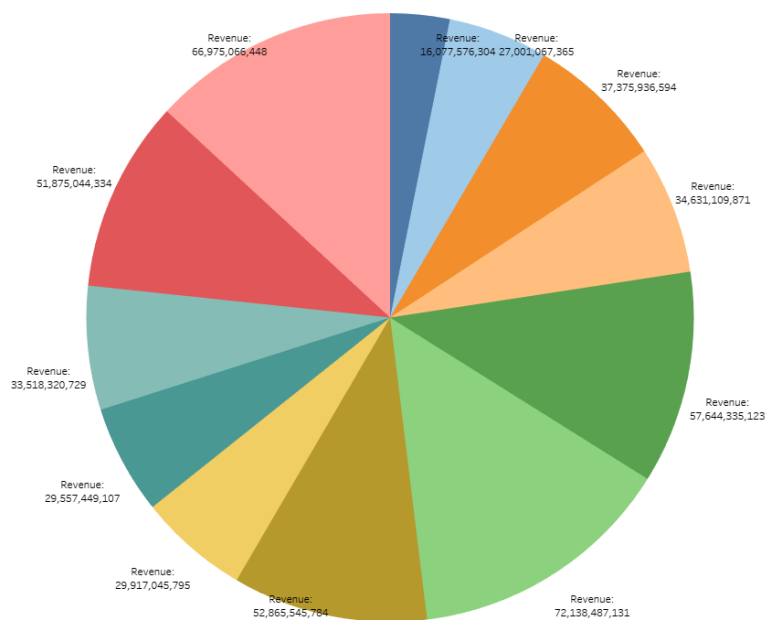


Top-10 Movie and Revenue





Months with maximum Revenue



SUM(Revenue)

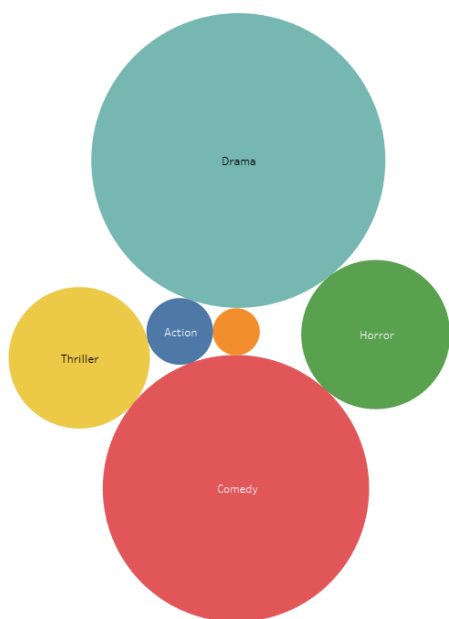
509,576,984,585

Month

January
February
March
April
May
June
July
August
September
October
November
December

Sheet 1 Sheet 2 Sheet 3 Sheet 4 Sheet 5 **Sheet 7**

Genre vs Movie Title



Genre

Action
Adventure
Comedy
Drama
Horror
Thriller

Sheet 1 Sheet 2 Sheet 3 Sheet 4 Sheet 5 Sheet 7 **Sheet 9**