# Book Project

## Create and Load data in tables

### Create Book Ratings table and Load data into it

```
create table if not exists bookratings
(userid string, isbn string, bookrating string)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
 "separatorChar" = "\;",
 "quoteChar" = '\"'
 ) LOCATION '/user/praveenkum79edu/Book_Analysis';

Load data inpath '/user/praveenkum79edu/Book_Data/BX-Book-Ratings.csv' INTO TABLE bookratings;
```



### Create Books Table and load data in to it

```
create table if not exists bookstable
(isbn string, title string, author string, year string, publisher string,urls
string,urlm string,urll string)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
 WITH SERDEPROPERTIES (
 "separatorChar" = "\;",
 "quoteChar" = '\"'
 )  LOCATION '/user/praveenkum79edu/Books_new_Analysis'
tblproperties ("skip.header.line.count"="1") ;

Load data inpath '/user/praveenkum79edu/Book_Data/BX-Books.csv' INTO TABLE bookstable;
```

# Problem Statement

✦ **Find out the frequency of books published each year. (Hint: Use Boooks.csv file for this)**

```
select year,count(*) as Published_Frequency from bookstable
group by year
order by cast(year as bigint);
```



✦ **Find out in which year maximum number of books were published**

```
select year, count(*) as A
from bookstable
group by year
order by A desc
limit 1
```



✦ **Find out how many book were published based on ranking in the year 2002. ( Hint: Use Book.csv and Book-Ratings.csv)**

```
select  bookrating, count(*) from bookstable
join bookratings on bookstable.isbn=bookratings.isbn
where  year=2002
group by bookratings.bookrating ;
```

Hive     maximum number of books were published     Add a description...

2m, 7s   Database default ▾   Type text ▾   ⚙   ?

```
1 select  bookrating, count(*) from bookstable
2 join bookratings on bookstable.isbn=bookratings.isbn
3 where   year=2002
4 group by bookratings.bookrating ;
5
6
```

```
INFO  : Number of reduce tasks not specified. Estimated from input data size: 3
INFO  : In order to change the average load for a reducer (in bytes):
INFO  :   set hive.exec.reducers.bytes.per.reducer=<number>
INFO  : In order to limit the maximum number of reducers:
INFO  :   set hive.exec.reducers.max=<number>
INFO  : In order to set a constant number of reducers:
INFO  :   set mapreduce.job.reduces=<number>
INFO  : number of splits:3
INFO  : Submitting tokens for job: job_1635139249191_4643
INFO  : Executing with tokens: []
INFO  : The url to track the job: http://ip-10-1-1-204.ap-south-1.compute.internal:8066/proxy/application_1635139249191_4643/
INFO  : Starting Job = job_1635139249191_4643, Tracking URL = http://ip-10-1-1-204.ap-south-1.compute.internal:8066/proxy/application_1635139249191_4643/
INFO  : Kill Command = /opt/cloudera/parcels/CDH-6.2.1-1.cdh6.2.1.p0.1425774/lib/hadoop/bin/hadoop job  -kill job_1635139249191_4643
```

job_1635139249191_4643
job_1635139249191_4644
job_1635139249191_4645

Query History     Saved Queries     Results (11)

| | | bookrating | _c1 |
|---|---|---|---|
| | 1 | 0 | 107630 |
| | 2 | 1 | 286 |
| | 3 | 10 | 12546 |
| | 4 | 2 | 520 |
| | 5 | 3 | 1068 |
| | 6 | 4 | 1722 |
| | 7 | 5 | 7228 |
| | 8 | 6 | 6372 |
| | 9 | 7 | 13340 |
| | 10 | 8 | 19762 |
| | 11 | 9 | 13128 |