```
o x
hadoop@0e8a607bbfab: / X proot@0e8a607bbfab: /home/h X + V
hadoop@0e8a607bbfab:/$ hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/hadoop/hive/lib/log4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/hadoop/hadoop/share/hadoop/common/lib/slf4j-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.apache.logging.slf4j.Log4jLoggerFactory]
Logging initialized using configuration in jar:file:/home/hadoop/hive/lib/hive-common-2.3.9.jar!/hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using H
ive 1.X releases.
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.hive.common.StringInternUtils (file:/home/hadoop/hive/lib/hive-common-2.3.9.jar) to field java.net.U
RI.string
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.hive.common.StringInternUtils
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
hive> SHOW DATABASES;
OK
default
tripdata
Time taken: 0.383 seconds, Fetched: 2 row(s)
hive> CREATE DATABASE IF NOT EXISTS tripdata_andres_leal_db;
Time taken: 0.057 seconds
hive> USE tripdata andres leal db;
Time taken: 0.024 seconds
hive>
```

```
hive> CREATE TABLE IF NOT EXISTS tripdata andres leal db.payments(vendorID INT, tpep pickup datetime TIMESTAMP, payment type STRING, total amount DOUBLE)
   > COMMENT 'PAYMENTS TABLE'
   > ROW FORMAT DELIMITED
   > FIELDS TERMINATED BY ',';
Time taken: 0.098 seconds
hive> CREATE TABLE IF NOT EXISTS tripdata andres leal db.passengers(tpep pickup datetime TIMESTAMP, passenger count INT, total amount DOUBLE)
   > COMMENT 'PASSENGERS TABLE'
   > ROW FORMAT DELIMITED
   > FIELDS TERMINATED BY ',';
Time taken: 0.056 seconds
hive> CREATE TABLE IF NOT EXISTS tripdata andres leal db.tolls(tpep pickup datetime TIMESTAMP, passenger count INT, tolls amount DOUBLE, total amount DOUBLE)
   > COMMENT 'TOLLS TABLE'
   > ROW FORMAT DELIMITED
   > FIELDS TERMINATED BY ',';
Time taken: 0.062 seconds
hive> CREATE TABLE IF NOT EXISTS tripdata andres leal db.congestion(tpep pickup datetime TIMESTAMP, passenger count INT, congestion surcharge DOUBLE, total amount DOUBLE)
   > COMMENT 'CONGESTION TABLE'
   > ROW FORMAT DELIMITED
   > FIELDS TERMINATED BY ',';
Time taken: 0.056 seconds
hive> CREATE TABLE IF NOT EXISTS tripdata andres leal db.distance(tpep pickup datetime TIMESTAMP, passenger count INT, trip distance DOUBLE, total amount DOUBLE)
   > COMMENT 'DISTANCE TABLE'
   > ROW FORMAT DELIMITED
   > FIELDS TERMINATED BY ',';
Time taken: 0.053 seconds
hive>
```

hadoop@0e8a607bbfab: / X hadoop@0e8a607bbfab: ~/lar X + V

```
hive> DESCRIBE FORMATTED passengers;
# col name
                        data_type
                                                 comment
tpep pickup datetetime timestamp
                        int
passenger count
                        double
total amount
# Detailed Table Information
Database:
                        tripdata andres leal db
                        hadoop
Owner:
CreateTime:
                        Sun Dec 24 16:10:30 ART 2023
LastAccessTime:
                        UNKNOWN
Retention:
Location:
                        hdfs://172.17.0.2:9000/user/hive/warehouse/tripdata andres leal db.db/passengers
Table Type:
                        MANAGED_TABLE
Table Parameters:
                                {\"BASIC_STATS\":\"true\"}
        COLUMN STATS ACCURATE
        comment
                                PASSENGERS TABLE
       numFiles
                                0
       numRows
                                0
        rawDataSize
                                0
        totalSize
                                0
        transient lastDdlTime
                                1703445030
# Storage Information
SerDe Library:
                        org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
                        org.apache.hadoop.mapred.TextInputFormat
InputFormat:
                        org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
OutputFormat:
Compressed:
                        No
Num Buckets:
                        -1
Bucket Columns:
                        Sort Columns:
                        []
Storage Desc Params:
        field.delim
        serialization.format
Time taken: 0.073 seconds, Fetched: 34 row(s)
hive>
```

hadoop@0e8a607bbfab: /

root@0e8a607bbfab: /home/h X

```
hadoop@0e8a607bbfab: /
                         root@0e8a607bbfab: /home/h X
hive> DESCRIBE FORMATTED distance:
OK
# col name
                        data type
                                                 comment
tpep pickup datetetime timestamp
passenger count
                        int
trip distance
                        double
total amount
                        double
# Detailed Table Information
Database:
                        tripdata andres leal db
                        hadoop
Owner:
CreateTime:
                        Sun Dec 24 16:10:50 ART 2023
LastAccessTime:
                        UNKNOWN
Retention:
Location:
                        hdfs://172.17.0.2:9000/user/hive/warehouse/tripdata andres leal db.db/distance
Table Type:
                        MANAGED TABLE
Table Parameters:
                                 {\"BASIC STATS\":\"true\"}
        COLUMN STATS ACCURATE
        comment
                                DISTANCE TABLE
        numFiles
                                 0
        numRows
        rawDataSize
                                 0
        totalSize
        transient lastDdlTime
                                 1703445050
# Storage Information
SerDe Library:
                        org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
                        org.apache.hadoop.mapred.TextInputFormat
InputFormat:
OutputFormat:
                        org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
Compressed:
                        No
Num Buckets:
                        -1
                        Bucket Columns:
Sort Columns:
                        []
Storage Desc Params:
        field.delim
        serialization.format
Time taken: 0.047 seconds, Fetched: 35 row(s)
hive>
```

```
hadoop@0e8a607bbfab: ~/la X
   hadoop@0e8a607bbfab: /
hadoop@0e8a607bbfab:~/landing$ hdfs dfs -ls /user/hadoop
Found 2 items
drwxr-xr-x

    hadoop supergroup

                                           0 2023-12-24 20:32 /user/hadoop/.sparkStaging
drwxr-xr-x - hadoop supergroup
                                           0 2022-01-22 23:56 /user/hadoop/inputs
hadoop@0e8a607bbfab:~/landing$ hdfs dfs -put yellow tripdata 2021-01.csv /user/hadoop
hadoop@0e8a607bbfab:~/landing$ hdfs dfs -ls /user/hadoop
Found 3 items
drwxr-xr-x

    hadoop supergroup

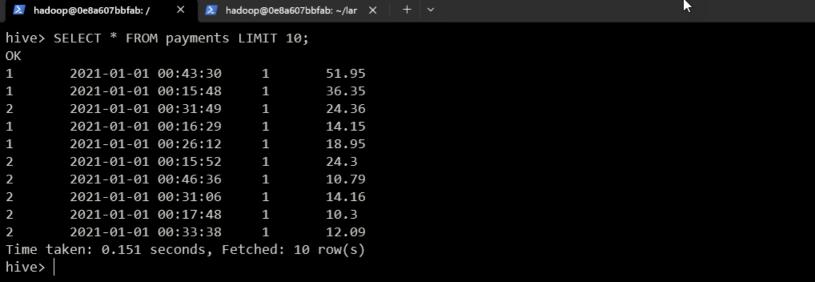
                                           0 2023-12-24 20:32 /user/hadoop/.sparkStaging
drwxr-xr-x - hadoop supergroup
                                           0 2022-01-22 23:56 /user/hadoop/inputs
                                  125981363 2023-12-24 20:33 /user/hadoop/yellow tripdata 2021-01.csv
            1 hadoop supergroup
-rw-r--r--
hadoop@0e8a607bbfab:~/landing$
```

```
X 💹 hadoop@0e8a607bbfab: ~/la X 🕂 🔻
hadoop@0e8a607bbfab: /
>>> df = spark.read.options(delimiter=";", header=True).csv("yellow tripdata 2021-01.csv")
>>> df.printSchema()
 -- VendorID, tpep pickup datetime, tpep dropoff datetime, passenger count, trip distance, RatecodeID, store and fwd flag, PULocationID, DOLocationID, payment type, fare amount, extra, mta tax, tip am
ount, tolls amount, improvement surcharge, total amount, congestion surcharge: string (nullable = true)
```

```
× Nadoop@0e8a607bbfab: ~/la ×
hadoop@0e8a607bbfab: /
>>> df = spark.read.options(header=True).csv("yellow tripdata 2021-01.csv")
>>> df.printSchema()
root
 |-- VendorID: string (nullable = true)
 -- tpep pickup datetime: string (nullable = true)
  -- tpep dropoff datetime: string (nullable = true)
  -- passenger count: string (nullable = true)
  -- trip distance: string (nullable = true)
 -- RatecodeID: string (nullable = true)
  -- store and fwd flag: string (nullable = true)
  -- PULocationID: string (nullable = true)
 -- DOLocationID: string (nullable = true)
  -- payment type: string (nullable = true)
  -- fare amount: string (nullable = true)
  -- extra: string (nullable = true)
 -- mta_tax: string (nullable = true)
 -- tip amount: string (nullable = true)
  -- tolls amount: string (nullable = true)
  -- improvement surcharge: string (nullable = true)
 |-- total amount: string (nullable = true)
 |-- congestion surcharge: string (nullable = true)
>>> df2 = spark.read.options(header=True, inferSchema=True).csv("yellow tripdata 2021-01.csv")
>>> df2.printSchema()
root
 -- VendorID: integer (nullable = true)
 |-- tpep pickup datetime: string (nullable = true)
  -- tpep dropoff datetime: string (nullable = true)
  -- passenger count: integer (nullable = true)
  -- trip distance: double (nullable = true)
 -- RatecodeID: integer (nullable = true)
  -- store and fwd flag: string (nullable = true)
  -- PULocationID: integer (nullable = true)
  -- DOLocationID: integer (nullable = true)
  -- payment_type: integer (nullable = true)
  -- fare amount: double (nullable = true)
  -- extra: double (nullable = true)
 -- mta tax: double (nullable = true)
  -- tip amount: double (nullable = true)
  -- tolls amount: double (nullable = true)
  -- improvement surcharge: double (nullable = true)
 -- total amount: double (nullable = true)
 -- congestion surcharge: double (nullable = true)
```

K

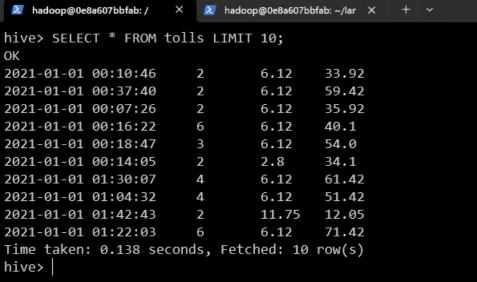
```
Madoop@0e8a607bsfab: ~/la ×
hadoop@0e8a607bbfab: /
>>> payments df credit card = payments df.filter(col("payment type")==1)
>>> payments df credit card.write.insertInto("tripdata andres leal db.payments")
2023-12-24 21:44:22,198 WARN conf. HiveConf: HiveConf of name hive.metastore.local does not exist
2023-12-24 21:44:22,656 WARN session. Session State: METASTORE FILTER HOOK will be ignored, since hive security authorization manager is set to instance of HiveAuthorizerFactory.
>>> payments df credit card.show(10)
  |VendorID|tpep pickup datetime|payment type|total amount|
       1 2021-01-01 00:43:30
                                                51.95
         2021-01-01 00:15:48
                                                36.35
       2 2021-01-01 00:31:49
                                                24.36
         2021-01-01 00:16:29
                                                14.15
         2021-01-01 00:26:12
                                                18.95
       2 2021-01-01 00:15:52
                                                24.3
       2 2021-01-01 00:46:36
                                                10.79
       2 2021-01-01 00:31:06
                                                14.16
         2021-01-01 00:17:48
                                                10.3
       2 2021-01-01 00:33:38
                                                12.09
only showing top 10 rows
>>>
```



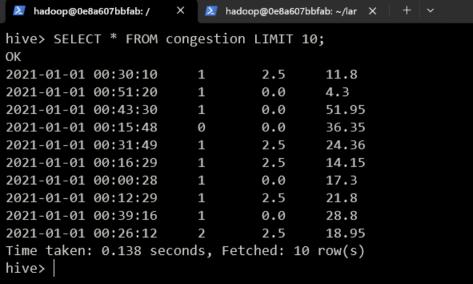
```
hadoop@0e8a607bbfab: ~/la X
 hadoop@0e8a607bbfab: /
>>> passengers df condition = passengers df.filter( ( col("passenger count") > 2 ) & ( col("total amount") > 8 ) )
>>> passengers df condition.write.insertInto("tripdata andres leal db.passengers")
>>> passengers df condition.show(10)
|tpep pickup datetime|passenger count|total amount|
 2021-01-01 00:15:52
                                              24.3
 2021-01-01 00:31:06
                                              14.16
 2021-01-01 00:42:11
                                               8.3
 2021-01-01 00:43:41
                                               9.3
 2021-01-01 00:34:37
                                               18.3
 2021-01-01 00:06:08
                                               13.3
 2021-01-01 00:19:57
                                              40.3
 2021-01-01 00:28:07
                                               14.8
 2021-01-01 00:08:04
                                              18.59
 2021-01-01 00:22:02
                                             13.56
only showing top 10 rows
```

hadoop@0e8a607bbfab: /	×	hadoop@0e8a607bbfab: ~/lar	x + ~
hive> SELECT * FROM	passeng	ers LIMIT 10;	
OK			
2021-01-01 00:15:52	3	24.3	
2021-01-01 00:31:06	5	14.16	
2021-01-01 00:42:11	5	8.3	
2021-01-01 00:43:41	3	9.3	
2021-01-01 00:34:37	4	18.3	
2021-01-01 00:06:08	4	13.3	
2021-01-01 00:19:57	3	40.3	
2021-01-01 00:28:07	5	14.8	
2021-01-01 00:08:04	3	18.59	
2021-01-01 00:22:02	3	13.56	
Time taken: 0.174 se	conds,	Fetched: 10 row(s))
hive>			

```
hadoop@0e8a607bbfab: /
                     X № hadoop@0e8a607bbfab: ~/la X
>>> tolls df condition = tolls df.filter( (col("tolls amount") > 0.1) & (col("passenger count") > 1) )
>>> tolls df condition.show(10)
|tpep pickup datetime|passenger count|tolls amount|total amount|
 2021-01-01 00:10:46
                                              6.12
                                                          33.92
 2021-01-01 00:37:40
                                              6.12
                                                          59.42
 2021-01-01 00:07:26
                                              6.12
                                                          35.92
 2021-01-01 00:16:22
                                                           40.1
                                              6.12
 2021-01-01 00:18:47
                                              6.12
                                                           54.0
 2021-01-01 00:14:05
                                               2.8
                                                           34.1
 2021-01-01 01:30:07
                                              6.12
                                                          61.42
 2021-01-01 01:04:32
                                              6.12
                                                          51.42
 2021-01-01 01:42:43
                                             11.75
                                                          12.05
 2021-01-01 01:22:03
                                              6.12
                                                          71.42
only showing top 10 rows
>>> tolls df condition.write.insertInto("tripdata andres leal db.tolls")
>>>
```



```
>>> congestion df condition = congestion df.filter( (col("tpep pickup datetime") > '2021-01-01 00:00') & (col("tpep_pickup_datetime") < '2021-01-02 00:00') )
>>> congestion df condition.show(10)
|tpep_pickup_datetime|passenger count|congestion surcharge|total amount|
 2021-01-01 00:30:10
                                                      2.5
                                                                  11.8
 2021-01-01 00:51:20
                                                      0.0
                                                                   4.3
 2021-01-01 00:43:30
                                                      0.0
                                                                 51.95
 2021-01-01 00:15:48
                                                                 36.35
                                                      0.0
 2021-01-01 00:31:49
                                                      2.5
                                                                 24.36
 2021-01-01 00:16:29
                                                      2.5
                                                                 14.15
 2021-01-01 00:00:28
                                                      0.0
                                                                  17.3
 2021-01-01 00:12:29
                                                      2.5
                                                                  21.8
 2021-01-01 00:39:16
                                                      0.0
                                                                  28.8
 2021-01-01 00:26:12
                                                      2.5
                                                                 18.95
only showing top 10 rows
>>> congestion df condition.write.insertInto("tripdata andres leal db.congestion")
>>>
```



```
≥ hadoop@0e8a607bbfab: / × ≥ hadoop@0e8a607bbfab: ~/lar × + ×
hive> SELECT * FROM distance;
OK
2020-12-31 21:40:20 1 17.96 53.3
Time taken: 0.144 seconds, Fetched: 1 row(s)
hive>
```