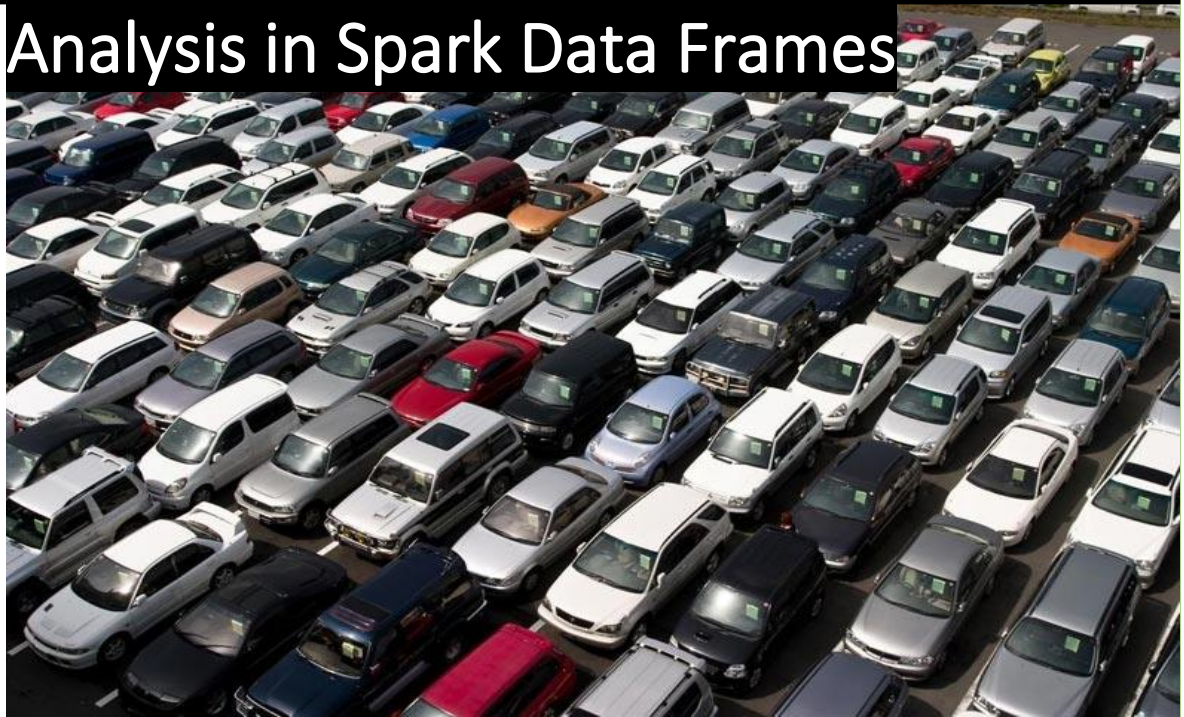# 2022

# Business Plan: Used Cars Dealership

## Analysis in Spark Data Frames



**VISHAL KAATAL**

3/14/2022

# TABLE OF CONTENTS

# TABLE OF FIGURES

# INTRODUCTION

The used car business is booming worldwide because people invest their money in properties, stocks, businesses instead of brand-new cars. (DriveNation, 2019) The brand-new vehicle loses over 10% of its value during the first month of purchase and up to 20% within the first year; therefore, people are leaning towards used cars to save these losses. A market survey has been conducted by Miroslav Zoricak and published the dataset called cars.csv on the Kaggle website under the name as classified ads for cars (Zoricak, 2017) which is the basis of this analysis report for setting up a used car dealership in the Czech Republic and Germany.

This report also covers cleaning the datasets based on the Business Plan requirement. The visualization graphs have been provided to understand the data to make an investment decision. Some of the analysis questions that have been discussed in the report are as follows:

1. How many cars are available in each vehicle manufacturing year between 2000 – 2017?

2. What is the top 10 highest average price of a vehicle based on make and model?

3. What is the top 10 lowest average price of a car based on make and model?

4. Describe and identify the top 5 make and models based on the average price of a car for three different segments as Economic, Intermediate and Luxury?

5. What kind of transmission do people prefer in the Czech Republic and Germany?

6. What are the top 5 vehicle manufacturers recommended to invest in? Why?

# DATA CLEANING

(Zoricak, 2017) The collected data has been scraped over one year, and sources are completely unstructured, so as a result, the information is missing values, and some values are wrong (E.g., phone numbers scraped as mileage or price). Therefore, the dataset must be cleaned to answer the analysis questions described earlier. The dataset consists of roughly 3.5 million rows and 16 columns. The 16 columns consist of entities such as maker, model, mileage, manufacture year, engine displacement, engine power, body type, colour type, last emission check year, transmission, door count, seat count, fuel type, data collection date, car last seen on the website, and price in euro. Out of 16 columns, only 13 answered the business plan in the cleaning and analysis because 50% or more values were missing in those removed columns. The cleaning steps have been discussed in the Appendix with details.

To start a used car business, selecting the famous car manufacturers (makers) is essential such that the customers are attracted to your dealership and ultimately make their purchase. Therefore, the manufacturing year for the car has to be between 2000 and 2017 included with a price range from €3000 to €2,000,000. Furthermore, it is divided into three segments. First is Economic segment customers, who would like to spend between €3000 - €20,000. Second is Intermediate segment customers, who have a budget for a car between €20,000 – €300,000. Finally, the Luxury segment customers can spend between €300,000 – €2,000,000 for a luxury vehicle. Hence, some people might prefer specific transmission or number of doors in a car. Therefore, Using these criteria, customers will be happy to select their dream car based on their vehicle's budget, requirement, and purpose, like a daily commute or high-end luxury cars for recreational purposes, etc.

# ANALYSIS RESULT

This part of the report answers the earlier analysis questions to make an investment decision.
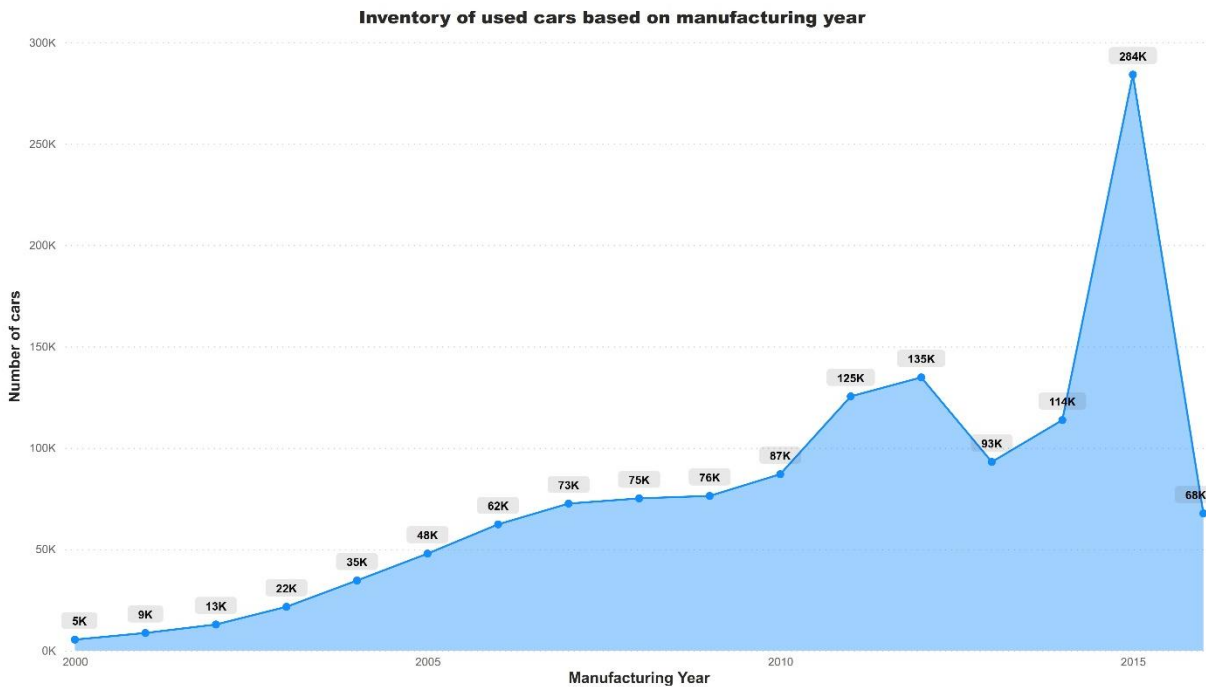


*Figure 1: Inventory of used cars*

Based on Figure 1, it can be stated that most used cars are available for the manufacturing year 2015. If anyone is looking for a newer used car, they should look for 2010 – 2016 maker cars with low mileage, better transmission quality, and low price.

Figure 2 represents the Top 10 average prices of a car based on maker and model, out of which Lamborgini Aventador has the highest average price of € 365,961. Lamborgini, Porche, BMW, Tesla, Bentley, and Rolls-Royce are the most expensive makers. On the other hand, figure 3 represents the Top 10 lowest average car prices. The least expensive car is the Skoda Galaxy, based on the analysis with an average value of €3071.80, whereas the most miniature makers are Fiat, Opel, Skoda. However, the least expensive ones are the most popular makers such as Volkswagen, Audi, Opel, Ford, Skoda and have the highest inventory count, as shown in Figure 4.

**Top 10 average price of car based on maker and model**

*Figure 2: Top 10 highest average prices of a car*



**Top 10 lowest average price of car based on maker**

*Figure 3: Top 10 lowest average prices of car*

*Figure 4: Top 10 most popular car makers*



*Figure 5: Top 5 Economic segment cars*

The average price can be from € 3000 - € 2,000,000 for the customers based on their needs and budget. It

has been divided into three segments of this analysis: Economic (€ 3000 - € 20,000), Intermediate (€ 20,000

- € 300,000) and Luxury (€ 300,000 - € 2,000,000). Figure 5 represents the Top 10 Economic segments cars

with their maker and model: Volvo V40, Peugeot RCZ, Skoda 130, Lexus GS, and Volkswagen Beetle. Figure 6 shows the Top 5 Intermediate segment cars with their maker and model: BMW Z8, Tesla Roadster, Tesla Model-X, Bentley Brooklands and Rolls-Royce Wraith. In the last, the Luxury segment has only two cars: Lamborgini Aventor and Porsche Carrera-GT, as shown in figure 7.



*Figure 6: Top 5 Intermediate segment cars*



*Figure 7: Luxury Segment cars*

From Figure 8, it can be derived that many cars in the Czech Republic and Germany are manual cars instead of automatic. Also, only two significant carmakers, Audi and Volkswagen, produce most cars with automatic transmission, whereas others tend to make manual cars. The reason behind manual vs automatic is the cost (price); as for all carmakers, cars with automatic transmission are more expensive than those with manual transmission. Also, the maximum number of cars produced in Manual transmission are Volkswagen, Opel, Ford and Skoda from these manufacturers.



*Figure 8: Top 10 car makers transmission type*

Lastly, in figure 9, we notice that most manufacturers produce cars with a maximum of 5 seats, equivalent to more than 75 % of the market.



*Figure 9: Number of seats preferred in a car*

# CONCLUSION

In general, any business plan has risks, whether its business model will be successful or not. During the analysis, it is clear that if the desired selection of used car maker, model, and manufacturing year is selected, the business shall not fail. All the parameters have been discussed in previous chapters, and appropriate graphs have been generated such that visualization should give us a better understanding. The most crucial part of data cleaning has been done using Spark such that the decisions can be data-driven for the investment in the operated car dealership business.

It can be concluded that if the investment firm is investing in this business plan, it can generate a considerable profit quickly. Therefore, to make it successful, the car manufacturers selected are Audi, Skoda, Ford, Seat, BMW, and Opel manufactured between 2000 – 2016 with a number of car seats above 4 with any transmission type. However, to offer some high-budget buyers, the dealership can include manufacturers such as Lamborghini, Porsche, Tesla, Bently, and Rolls-Royce to attract customers into the dealership. Based on the sales pitch, deals can be secured to earn more profit from Intermediate and luxury cars. The best practice is to include variety and price range so customers come to the dealership. The salesperson can secure them a deal that will make them happy and comfortable to buy again in the future.

# BIBLIOGRAPHY

*DriveNation*. (2019, 05 03). Retrieved from drivenation.ca/blog:

https://www.drivenation.ca/blog/depreciation-vehicle-much-money-will-lose-buying-new-car/

Zoricak, M. (2017, 03 16). *Classified Ads for Cars*. Retrieved from Kaggle:

https://www.kaggle.com/mirosval/personal-cars-classifieds

# APPENDIX

////Loading of Data in GCP///

wget https://www.dropbox.com/s/rsrxro7r1c5a4i2/cars.csv

////Moving Data into HDFS////

hadoop fs -mkdir /BigData

hadoop fs -copyFromLocal cars.csv /BigData/.

////Analysis in Spark DataFrame////

spark-shell --master yarn

/* Run import command for sql functions */

```
import org.apache.spark.sql.functions.{expr, col, column, min, max,
desc, avg, sum}
```

```
import org.apache.spark.sql.types._
```



*Task 1: Write a Spark DataFrames query to create a table called **used_cars** from data. Use a schema that is appropriate for the column headings*

/* Loading of a cars dataset using schema in Spark as */

-- use :paste for the following set of instructions

```
val schema = StructType(Array(
        StructField("maker", StringType, true),
        StructField("model", StringType, true),
        StructField("mileage", IntegerType, true),
        StructField("manufacture_year", IntegerType, true),
        StructField("engine_displacement", IntegerType, true),
        StructField("engine_power", IntegerType, true),
```

```
            StructField("body_type", StringType, true),
            StructField("color_slug", StringType, true),
            StructField("stk_year", IntegerType, true),
            StructField("transmission", StringType, true),
            StructField("door_count", IntegerType, true),
            StructField("seat_count", IntegerType, true),
            StructField("fuel_type", StringType, true),
            StructField("date_created", DateType, true),
            StructField("datelastseen", DateType, true),
            StructField("price_eur", FloatType, true)))

val used_cars = spark
  .read.format("csv")
  .option("header", "True")
  .schema(schema)
  .load("hdfs://10.128.0.6:8020/BigData/cars.csv")
```

-- end paste (Cntrl-D)

```
used_cars.show
```

*Task 2: Look at the date column of the table used_cars. Why does the date column have all NULL or incomplete values?*

The data type selected for a date is DateType, due to which in-completed values of date have been shown as output. Even though it still shows the proper dates instead of nulls, it can be fixed in the next task.

*Task 3: Create a table such that the date column is read correctly based on the format in the dataset.*

```
/* Fixing the DataType for dates and creating new schema. Loading the dataset again into spark */
-- use :paste
val schema_new = StructType(Array(
            StructField("maker", StringType, true),
            StructField("model", StringType, true),
            StructField("mileage", IntegerType, true),
            StructField("manufacture_year", IntegerType, true),
            StructField("engine_displacement", IntegerType, true),
            StructField("engine_power", IntegerType, true),
            StructField("body_type", StringType, true),
            StructField("color_slug", StringType, true),
            StructField("stk_year", IntegerType, true),
            StructField("transmission", StringType, true),
            StructField("door_count", IntegerType, true),
            StructField("seat_count", IntegerType, true),
            StructField("fuel_type", StringType, true),
            StructField("date_created", TimestampType, true),
            StructField("datelastseen", TimestampType, true),
            StructField("price_eur", FloatType, true)))

val used_cars = spark
  .read.format("csv")
  .option("header", "True")
  .schema(schema_new)
  .load("hdfs://10.128.0.6:8020/BigData/cars.csv")
-- end paste (Cntrl-D)

used_cars.show
```



In the above figure, the datatype for a date has been solved, which displays the date and time.

*Task 4: Write Spark DataFrames queries to see how many missing values you have in each attribute?
Based on the results, document how many missing values in each column we have. Especially, mention
those columns with more than 50% missing values.*

/* Finding number of null values in dataset */
```
val number_of_null_values = used_cars.select(used_cars.columns.map(c =>
sum(col(c).isNull.cast("int")).alias(c)): _*).show
```



The current dataset contains 3.5 million records in each column; therefore, anything over 1.7 million null
values columns represents more than 50% missing values that are colour_slug (3.34 million), stk_year
(3.01 million) and fuel_type (1.84 million).

*Task 5: Group the price column and count the number of unique prices. Do you notice if there is a single
price repeating across the ads?*

/* Group the price column and count the number of unique prices */
-- use :paste for the following set of instructions
```
val cars_unique_prices =
used_cars.groupBy("price_eur").count.filter("count > 1")
val cars_unique_prices_desc =
cars_unique_prices.orderBy(col("count").cast("int").desc)
cars_unique_prices_desc.show
```
-- end paste (Cntrl-D)



The price of 1295.34 has a maximum count of 673623.

*Task 6: Write a Spark DataFrames query to create a new table called* **clean_used_cars** *from* **used_cars** *with the following conditions:*

- *Drop the columns with more than 50% missing values*
- *The manufacture year between 2000 and 2017, including 2000 and 2017*
- *Both maker and model exist in the row*
- *The price range is from 3000 to 2000,000 (3000 ≤ price ≤ 2000,000)*
- *Remove any price you singled out in Task 5 (i.e. a price that repeats too frequently for a random set of ads).*

/* Removing the col with more than 50 % null values i.e. color_slug, stk_year, fuel_type */
```
val clean_used_cars1 =
used_cars.drop("color_slug","stk_year","fuel_type")
```

/* Manufacture Year between 2000 and 2017 (inclusive) */
```
val clean_used_cars2 = clean_used_cars1.filter("manufacture_year >= 2000
AND manufacture_year <=2017")
```

/* Both maker and model exist in the row */
```
val clean_used_cars3 = clean_used_cars2.filter(col("maker").isNotNull)
val clean_used_cars4 = clean_used_cars3.filter(col("model").isNotNull)
```

/* Price range for car is from 3000 to 2000,000 */
```
val clean_used_cars5 = clean_used_cars4.filter("price_eur >= 3000 AND
price_eur <=2000000")
```

/* Removing price_eur = 1295.34 in the Task #5 and removing this price from the dataset. If you notice the highest one 1295.34 got removed during previous filtering process. Just for filteration, perform the following query */
```
val clean_used_cars = clean_used_cars5.filter("price_eur != 1295.34")

clean_used_cars.show(5)
```

*Task 7: Write a Spark DataFrames query to find how many records remained **clean_used_cars***

/* How many records remaining in clean_used_cars dataframe */

```
val number_of_null_values_clean =
clean_used_cars.select(clean_used_cars.columns.map(c =>
sum(col(c).isNull.cast("int")).alias(c)): _*).show
```



```
val number_of_records = clean_used_cars.count
```



The number of records remaining in the clean_used_cars is **1322853**.

*Task 8: Write a Spark DataFrames query to find the make and model for the cars with the top 10 **highest average price***

/* Spark DataFrames query, find the make & model for the cars with the top 10 highest average price */

```
val top_10_highest_avg_price = clean_used_cars
.select("maker","model","price_eur").groupBy("maker","model")
.agg(avg(col("price_eur")).alias("avg_price_eur"))
.orderBy(col("avg_price_eur").desc)
top_10_highest_avg_price.show(10)
```

*Task 9: Write a Spark DataFrames query to find the make and model for the cars with the top 10 **lowest average price***

```
/* Spark DataFrames query, find the make & model for the cars with the top 10 lowest average price */
val top_10_lowest_avg_price = clean_used_cars
.select("maker","model","price_eur").groupBy("maker","model")
.agg(avg(col("price_eur")).alias("avg_price_eur"))
.orderBy(col("avg_price_eur").asc)

top_10_lowest_avg_price.show(10)
```



*Task 10 – 12: Write a Spark DataFrames query to recommend the top five make and models for different segments for customers*

- ***Economic** segment customers (3000≤price<20,000) - based on the top **average price***
- ***Intermediate** segment customers (20,000≤price<300,000) - based on the top **average price***
- ***Luxury** segment customers (300,000≤price<2000,000) - based on the top **average price***

```
/* Economic Segment DataFrame*/
val used_car_segment = clean_used_cars
.groupBy("maker","model")
.agg(avg(col("price_eur")).alias("avg_price_eur"))
.orderBy(col("avg_price_eur").desc)
```

```
/* Economic Segment Customers | Price: 3000 - 20000*/
val used_car_economic = used_car_segment.filter("avg_price_eur >=3000
AND avg_price_eur <20000")
used_car_economic.show(5)
```

/* Intermediate Segment Customers | Price: 20000 - 300000*/

```
val used_car_intermediate = used_car_segment.filter("avg_price_eur
>=20000 AND avg_price_eur <300000")
used_car_intermediate.show(5)
```

/* Luxury Segment Customers | Price: 300000 - 2000000*/
```
val used_car_luxury = used_car_segment.filter("avg_price_eur >=
300000")
used_car_luxury.show(5)
```

/* Inventory of used cars */

```
val manuf_count = clean_used_cars.groupBy("manufacture_year").count
val manuf_count_asc =
manuf_count.orderBy(col("manufacture_year").asc).show
```



/* Top 10 most popular car makers */

```
val popular_makers = clean_used_cars.groupBy("maker").count
val top10_popular_makers =
popular_makers.orderBy(col("count").desc).show(10)
```

```
scala> val popular_makers = clean_used_cars.groupBy("maker").count
popular_makers: org.apache.spark.sql.DataFrame = [maker: string, count: bigint]

scala> val top10_popular_makers = popular_makers.orderBy(col("count").desc).show(10)
+----------+------+
|     maker| count|
+----------+------+
|volkswagen|166585|
|      audi|143998|
|      opel|123454|
|      ford|113771|
|     skoda|100824|
|      fiat| 71159|
|   citroen| 68983|
|    renault| 47249|
|   peugeot| 44059|
|    toyota| 43470|
+----------+------+
only showing top 10 rows

top10_popular_makers: Unit = ()
```

/* Car makers transmission type*/
```
val transmission_type_clean =
clean_used_cars.filter(col("transmission").isNotNull)
val transmission_type =
transmission_type_clean.select("maker","transmission").groupBy("maker"
,"transmission").count
val transmission_type_desc =
transmission_type.orderBy(col("count").desc).show
```



```
scala> :paste
// Entering paste mode (ctrl-D to finish)

/* Car makers transmission type*/
val transmission_type_clean = clean_used_cars.filter(col("transmission").isNotNull)
val transmission_type = transmission_type_clean.select("maker","transmission").groupBy("maker","transmission").count
val transmission_type_desc = transmission_type.orderBy(col("count").desc).show

// Exiting paste mode, now interpreting.

+----------+------------+------+
|     maker|transmission| count|
+----------+------------+------+
|volkswagen|         man|117693|
|      opel|         man|103656|
|      ford|         man| 93887|
|     skoda|         man| 73829|
|      audi|         man| 68785|
|      audi|        auto| 64750|
|      fiat|         man| 63795|
|   citroen|         man| 50472|
|volkswagen|        auto| 41462|
|    renault|         man| 39658|
|   peugeot|         man| 35525|
|     nissan|         man| 34652|
|      seat|         man| 34277|
|    toyota|         man| 30882|
|       bmw|        auto| 27432|
|   hyundai|         man| 25281|
|      mini|         man| 23293|
|     mazda|         man| 16449|
|       kia|         man| 15956|
|     skoda|        auto| 15344|
+----------+------------+------+
only showing top 20 rows

transmission_type_clean: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [maker: string, model: string ... 11 more fields]
transmission_type: org.apache.spark.sql.DataFrame = [maker: string, transmission: string ... 1 more field]
transmission_type_desc: Unit = ()
```
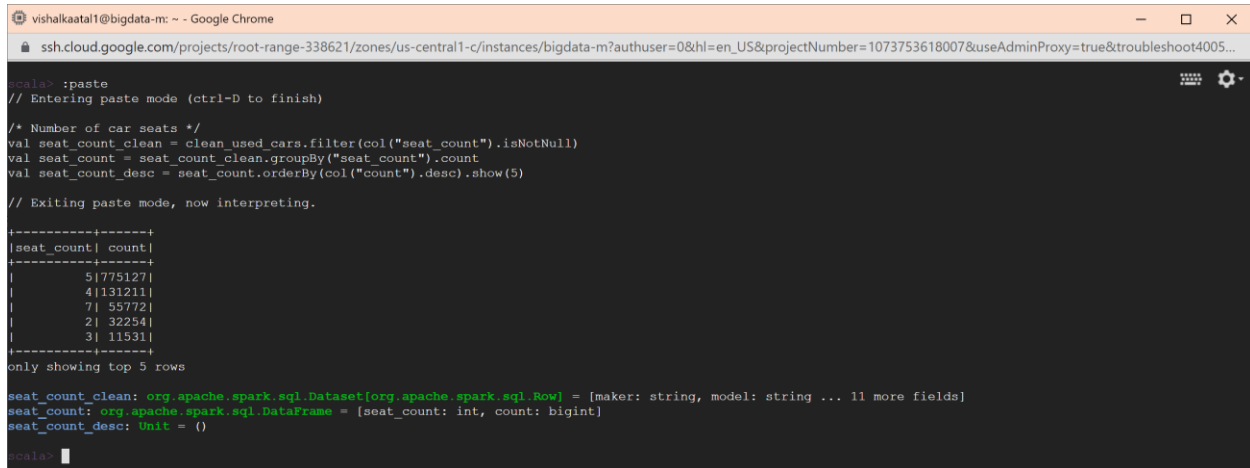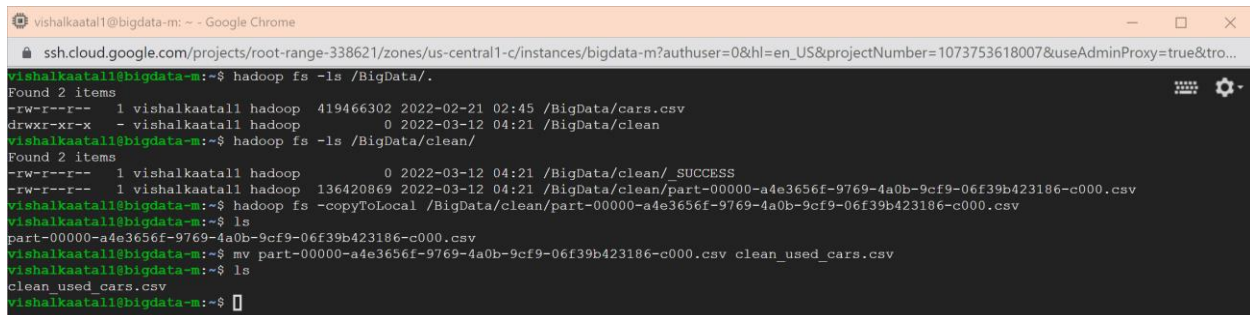
/* Number of car seats */

```
val seat_count_clean =
clean_used_cars.filter(col("seat_count").isNotNull)
val seat_count = seat_count_clean.groupBy("seat_count").count
val seat_count_desc = seat_count.orderBy(col("count").desc).show(5)
```



/* Download the dataset for visualization using the following commands */

```
clean_used_cars.coalesce(1).write
.option("header", true).csv("hdfs://10.128.0.6:8020/BigData/clean")
```



After moving the clean_used_cars dataset into GCP using the above commands, download the file directly to create all the visualizations in the Analysis section using Power BI.