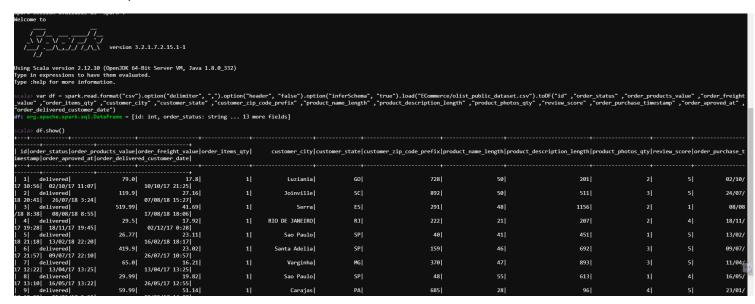
Loading dataset into FTP location and then to HDFS Directory:

```
[therohitsaha08gmail@sl-cdp-prod-en10 ~]$ cd Capstone/ECommerce
[therohitsaha08gmail@sl-cdp-prod-en10 ECommerce]$ ls
olist_public_dataset.csv
[therohitsaha08gmail@sl-cdp-prod-en10 ECommerce]$ hdfs dfs -ls
Found 6 items
drwxrwxr-x - therohitsaha08gmail hadoop
                                                0 2024-01-29 11:42 .sparkStaging
drwxrwxr-x - therohitsaha08gmail hadoop
                                                0 2024-01-25 08:41 ECommerce
drwxrwxr-x - therohitsaha08gmail hadoop
                                               0 2024-01-24 11:31 Weather_Analysis_Project
drwxrwxr-x - therohitsaha08gmail hadoop
                                               0 2024-01-18 12:12 YoutubeDataProject
drwxrwxrwx - therohitsaha08gmail hadoop
                                               0 2023-08-02 09:59 hive_demo
drwxrwxr-x - therohitsaha08gmail hadoop
                                               0 2023-04-22 13:18 simplilearn
[therohitsaha08gmail@sl-cdp-prod-en10 ECommerce]$ hdfs dfs -ls ECommerce
Found 1 items
            2 therohitsaha08gmail hadoop
                                          10304725 2024-01-25 08:41 ECommerce/olist public dataset.csv
rw-rw-r--
[therohitsaha08gmail@sl-cdp-prod-en10 ECommerce]$
```

Load dataset into Spark DataFrame from HDFS:



Extract only 'date' from all 3 timestamp columns respectively and store in a new dataframe:

		haseDate",to_date(unix_timestamp(col(pproveDate",to_date(unix_t	imestamp(col("order_	_aproved_at"),	"dd/MM/yy H:mm
		veryDate",to_date(unix_timestamp(col(" = [id: int, order_status: string		ustomer_date"), "de	d/MM/yy H:mm").cast("times	tamp")))				
scala> date_df.show()										
							+			
id order_status order_produ imestamp order_aproved_at ord	cts_value c er_delivere	order_freight_value order_items_qty ed_customer_date PurchaseDate ApproveD	customer_city date DeliveryDate				roduct_description_length			
		17.8 1			728	501	201	21	5	02/10/
17 10:56 02/10/17 11:07		10/10/17 21:25 2017-10-02 2017-10								
2 delivered L8 20:41 26/07/18 3:24	119.9	27.16 1 07/08/18 15:27 2018-07-24 2018-07	Joinville	sc	892	50	511	3	5	24/07/
3 delivered	519.99	41.69 1	2010-00-07 Serra	ES	291	48	1156	2	1	08/08
18 8:38 08/08/18 8:55		17/08/18 18:06 2018-08-08 2018-08								
4 delivered 7 19:28 18/11/17 19:45	29.5	17.92 1 02/12/17 0:28 2017-11-18 2017-11	RIO DE JANEIRO -18 2017-12-02	RJ	222	21	207	2	4	18/11/
5 delivered	26.77	23.11 1	Sao Paulo	SP	40	41	451	1	5	13/02/
8 21:18 13/02/18 22:20		16/02/18 18:17 2018-02-13 2018-02								
6 delivered 7 21:57 09/07/17 22:10	419.9	23.02 1 26/07/17 10:57 2017-07-09 2017-07	Santa Adelia Santa Adelia 30-26-	SP	159	46	692	3	5	09/07/
7 delivered	65.0	16.21 1	Varginha	MG	370	47	893	3	5	11/04/
17 12:22 13/04/17 13:25	20.001	13/04/17 13:25 2017-04-11 2017-04		col	401	ccl	canl	41	41	46 (05 (
8 delivered 7 13:10 16/05/17 13:22	29.99	19.82 1 26/05/17 12:55 2017-05-16 2017-05	Sao Paulo Sao Paulo Sao Paulo	SP	48	55	613	11	4	16/05/
9 delivered	59.99	51.14 1	Carajas	PA	685	28	96	4		23/01/
7 18:29 25/01/17 2:50	EC 001	02/02/17 14:08 2017-01-23 2017-01		l ra	2251	FOL	CALL	21	51	20.407.4
10 delivered 7 11:55 29/07/17 12:05	56.99	16.13 1 16/08/17 17:14 2017-07-29 2017-07	Resende Resende 2017-08-16	RJ	275	59	641	31	31	29/07/
11 delivered	599.0	15.69 1	Guaratingueta	SP	125	43	1686	2	4	16/05/

1. DAILY INSIGHTS

a. Sales

Total sales:

```
scala> var sales_df = date_df.groupBy(col("PurchaseDate")).agg(round(sum(col("order_products_value")),2).as("Total_Sales")).orderBy(asc("PurchaseDate"))
sales_df: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [PurchaseDate: date, Total_Sales: double]
  cala> sales_df.show()
|PurchaseDate|Total_Sales|
   2016-09-04
2016-09-05
                         200.0
   2016-09-13
                           99.0
   2016-09-15
                           90.0
                         56.99
962.3
   2016-10-02
2016-10-03
                        8827.6
   2016-10-04
   2016-10-05
                       6264.19
   2016-10-06
                       8214.07
   2016-10-07
                       6282.83
   2016-10-08
                       4575.48
   2016-10-09
                       2736.38
   2016-10-10
                       4130.77
   2016-10-22
                         49.0
                          82.5
   2016-12-23
    2017-01-05
                       3891.87
                        391.98|
357.36|
552.58|
   2017-01-06
2017-01-07
   2017-01-08
   2017-01-09
                        559.88
only showing top 20 rows
```

Total sales in each city:

```
cala> var sales_city_df = date_df.groupBy(col("PurchaseDate"),col("customer_city")).agg(round(sum(col("order_products_value")),2).as("Sales_per_City")).orderBy(asc("PurchaseDate")
ales_city_df: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [PurchaseDate: date, customer_city: string ... 1 more field]
   ala> sales_city_df.show()
                                    customer_city|Sales_per_City|
+----+-
|PurchaseDate|
    2016-09-04|
2016-09-05|
2016-09-13|
2016-09-15|
2016-10-02|
2016-10-03|
2016-10-03|
2016-10-03|
                                                                                        200.0
                                        Foz do Iguacu
Betim
                                  Betim
Feira de Santana
Fortaleza
Montes Claros
Eunapolis
Jundiai
                                                                                        99.0|
90.0|
56.99|
139.9|
45.8|
81.9|
58.9|
455.0|
42.0|
89.9|
48.9|
267.4|
330.0|
45.9|
18.9|
                                       Belo Horizonte
                                      Taubate
Taubate
Bauru
Sao Paulo
Estrela do Sul
Guacui
Recife
    2016-10-03|
2016-10-03|
2016-10-03|
2016-10-03|
    Canoas
Uberlandia
Campinas
                                                                                        46.99
30.9
34.9
     2016-10-04
2016-10-04
         showing top 20 rows
```

Total sales in each state:

b. Orders

Total number of orders:

```
var total_order = date_df.groupBy("PurchaseDate").agg(count("*").as("Total_Orders")).orderBy(asc("")).asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc("").asc
   otal_order: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [PurchaseDate: date, Total_Orders: bigint]
      cala> total_order.show()
|PurchaseDate|Total Orders|
            2016-09-04
                                                                                                                      1|
1|
1|
             2016-09-05
             2016-09-13
             2016-09-15
                                                                                                                      1
1
             2016-10-02
                                                                                                                 8
63
48
             2016-10-03
            2016-10-04
            2016-10-05
                                                                                                                  51
             2016-10-06
                                                                                                                  46
             2016-10-07
             2016-10-08
                                                                                                                  43
             2016-10-09
                                                                                                                  26
             2016-10-10
                                                                                                                   39
                                                                                                                 1|
1|
32|
4|
4|
6|
             2016-10-22
           2016-12-23
2017-01-05
            2017-01-06
             2017-01-07
             2017-01-08
             2017-01-09
                                                                                                                      5 j
only showing top 20 rows
```

City-wise order distribution:

```
var city_order = date_df.groupBy(col("PurchaseDate"),col("customer_city")).agg(count("*").as("Total_Orders")).orderBy(asc
rder: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [PurchaseDate: date, customer_city: string ... 1 more field
city_order
                                                                                                                                                                                                                          1 more field]
   ala> city_order.show()
                                    customer_city|Total_Orders|
PurchaseDate
   2016-09-04
2016-09-05;
2016-09-13;
2016-09-15;
2016-10-02;
2016-10-03;
2016-10-03;
2016-10-03;
2016-10-03;
2016-10-03;
2016-10-03;
2016-10-04;
2016-10-04;
                                    Foz do Iguacu
Betim
                                                                                  1|
1|
1|
1|
1|
1|
1|
1|
1|
1|
1|
1|
                              Feira de Santana
| Fortaleza
                                  Montes Claros
Belo Horizonte
                                              Guacui
                                  Eunapolis
Sao Paulo
Estrela do Sul
                                            Taubate
                                                  Bauru
                                               Jundiai
|
| Serra
|
| Garca
     2016-10-04
                                              Campinas
    2016-10-04
2016-10-04
                                             Anapolis
Buritis
     2016-10-04 Santana de Parnaiba
        showing top 20 ro
```

State-wise order distribution:

```
var state_order = date_df.groupBy(col("PurchaseDate"),col("customer_state")).agg(count("*").as("Total_Orders")).orderBy(asc(order: org.apache.spark.sql.@asc(order: org.apache.spark.sql.@asc(org.apache.spark.sql.@asc). 1 more field]
 tate_order: org.apache.spark.sql.
   ala> state order.show()
|PurchaseDate|customer state|Total Orders|
    2016-09-04
                                                                     1|
1|
1|
1|
1|
2|
4|
24|
4|
4|
1|
1|
1|
1|
   2016-09-04|

2016-09-05|

2016-09-13|

2016-09-15|

2016-10-02|

2016-10-03|
                                             MG
BA
                                             MG
ES
   2016-10-03|
2016-10-03|
2016-10-03|
2016-10-04|
                                             BA
MG
SP
SP
    2016-10-04
2016-10-04
                                             PE
RJ
                                             MG
RS
    2016-10-04
2016-10-04
    2016-10-04
2016-10-04
                                             GO
SC
                                             DF
PR
    2016-10-04
    2016-10-04
                                             ES|
BA|
    2016-10-04
    2016-10-04
 nly showing top 20 rows
```

Average review score per order:

```
var avg_rev = date_df.groupBy(col("PurchaseDate")).agg(round(sum(col("review_score"))/count("*"),3).as("Avg_Rev_per_Order")).orderBy(asc("PurchaseDate"))
v: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [PurchaseDate: date, Avg_Rev_per_Order: double]
 vg rev: org
        avg_rev.show()
|PurchaseDate|Avg_Rev_per_Order|
   2016-09-04
                                     3.0
   2016-09-05
                                    4.0
5.0
5.0
   2016-09-13
2016-09-15
   2016-10-02
   2016-10-03
2016-10-04
                                 3.625
3.857
   2016-10-05
   2016-10-06
2016-10-07
                                  4.0
3.739
   2016-10-08
                                  4.233
                                 3.923
   2016-10-09
   2016-10-10
                                  4.308
                                    1.0
   2016-10-22
   2016-12-23
   2017-01-05
                                  4.125
   2017-01-06
                                   3.25
   2017-01-07
                                    3.0
   2017-01-08
   2017-01-09
  ly showing top 20 rows
```

Average freight charges per order:

```
var avg_fre = date_df.groupBy(col("PurchaseDate")).agg(round(sum(col("order_freight_value"))/count("*"),3).as("Avg_Frei_per_Order")).orderBy(asc("PurchaseDate"))
a: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [PurchaseDate: date, Avg_Frei_per_Order: double]
 vg fre: or
  :ala> avg_fre.show()
|PurchaseDate|Avg_Frei_per_Order|
    2016-09-04
                                        39.88
14.16
    2016-09-05
   2016-09-13|
2016-09-15|
                                        19.93
26.89
   2016-10-02|
2016-10-03|
2016-10-04|
                                         15.15
                                      15.705
20.357
23.849
21.093
   2016-10-05
2016-10-06
    2016-10-07
                                       19.256
                                      24.621
23.502
    2016-10-08
    2016-10-09
   2016-10-10
2016-10-22
                                      21.252
37.97
   2016-12-23
2017-01-05
                                      44.16
20.032
                                      15.883
17.373
    2017-01-06
    2017-01-07
    2017-01-09
                                       16.624
  nly showing top 20 rows
```

Average time taken to approve the orders (order approved – order purchased):

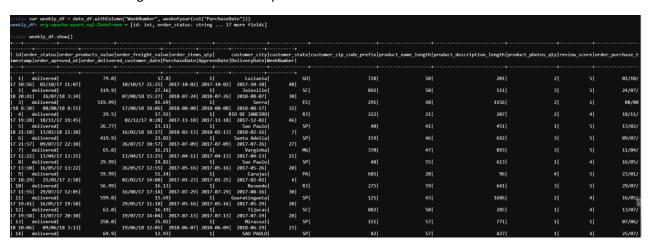
```
| The part of a graph of a graph
```

Average order delivery time:

```
### of delivery_diff = df withColour("coder_purchase_timestamp("), '64/MP/by Ham"), withColour("colour-delivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_colivered_col
```

2. WEEKLY INSIGHTS

new dataframe having conversion of 'order purchase date' to 'week number' :



a. Sales

Total sales:

```
cala> var weeklysales_df = weekly_df.groupBy("WeekNumber").agg(round(sum(col("order_products_value")),2).as("Total_Weekly_Sales")).orderBy("WeekNumber
eeklysales_df: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [WeekNumber: int, Total_Weekly_Sales: double]
  ala> weeklysales_df.show()
.
|WeekNumber|Total_Weekly_Sales|
                                 173944.85
             1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
                                 236810.92
255674.81
                                 265820.53
268544.45
                                  256302.52
                                  283705.55
                                  308138.64
                                   293062.7
                                  305988.97
                                  280830.14
                                  258832.68
                                 295689.26
293175.73
                                  308817.36
363995.98
 nly showing top 20 rows
```

Total sales in each city:

Total sales in each state:

b. Orders

Total number of orders:

```
scala> var totalweekly_order = weekly_df.groupBy("WeekNumber").agg(count("*").as("Total_Weekly_Orders")).orderBy("WeekNumber
totalweekly_order: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [WeekNumber: int, Total_Weekly_Orders: bigint]
  cala> totalweekly_order.show()
|WeekNumber|Total_Weekly_Orders|
                                     1439
             2
3
                                     1875
                                     1964
            4|
5|
6|
7|
8|
9|
                                     1956
                                     2077
                                     2148
                                     2060
                                     2145
                                     2412
            10
                                     2246
            11
                                     2221
                                     2381
            13
                                     2047
            14
                                     2190
            15
                                     2005
            16
                                     2258
            17
                                     2359
            18
                                     2570
            19
                                     2802
            201
                                     27971
only showing top 20 rows
```

City-wise order distribution:

```
cala> var weekly_city_order = weekly_df.groupBy(col("WeekNumber"),col("customer_city")).agg(count("*").as("Total_Weekly_Orders")).orderBy("WeekNumber"
eekly_city_order: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [WeekNumber: int, customer_city: string ... 1 more field]
  ala> weekly_city_order.show()
                         customer_city|Total_Weekly_Orders|
| WeekNumber |
                                                                       2
6
2
113
                               Americana
                  Balneario Camboriu
                         RIO DE JANEIRO
                                     Belem
                                                                          8 | 4 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 5 | 1 | 2 |
                             Santa Maria
                       Itaquaquecetuba
                  Sao Miguel Do Oeste
              1|
1|
                        Cachoeira Alta
                    Sao Jose do Egito
             1|Vitoria da Conquista|
1| Papucaia|
1| Fernando Prestes|
              1 |
1 |
                      Pontes e Lacerda
                                 Itaocara
              1 j
                            Rondonopolis
                  Bataguassu
| Sao Jose dos Campos
              1|
1|
              1|
1|
                           Ilha Solteira
                                   Cajamar
only showing top 20 rows
```

State-wise order distribution:

```
cala> var weekly_state_order = weekly_df.groupBy(col("WeekNumber"),col("customer_state")).agg(count("*").as("Total_Weekly_Orders")).orderBy("WeekNumber"
eekly_state_order: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [WeekNumber: int, customer_state: string ... 1 more field]
 cala> weekly_state_order.show()
| WeekNumber|customer_state|Total_Weekly_Orders
                1|
1|
1|
                                          TOI
                                                                               5 |
2 |
42 |
6 |
6 |
3 |
28 |
78 |
20 |
                                          SE
                                          PI
RO
                 1
1
1
                                          RN
                                          AM
                 1
1
1
1
1
1
1
1
1
                                          PE |
PR |
MT |
MS |
MG |
AC |
SP |
ES |
PB |
CE |
                                                                              20
170
2
                                                                              592
                                                                               7|
19|
                                                                              195
                                          PA BA
                1|
1|
                                                                               13
                                                                               52
 nly showing top 20 rows
```

Average review score per order:

```
cala> var weekly_avg_rev = weekly_df.groupBy(col("WeekNumber")).agg(round(sum(col("review_score"))/count("*"),3).as("Avg_Rev_per_Order")).orderBy("WeekNumber
eekly_avg_rev: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [WeekNumber: int, Avg_Rev_per_Order: double]
         weekly_avg_rev.show()
WeekNumber Avg_Rev_per_Order
                                      4.054
               1|
2|
3|
4|
5|
6|
7|
8|
9|
                                      4.06
4.03
                                     4.076
4.09
4.041
                                     4.057
4.042
                                      4.053
             10
11
12
13
14
15
16
17
18
19
                                     4.054
4.059
                                     4.071
                                      4.039
                                      4.029
                                      4.029
                                      4.105
                                      4.043
only showing top 20 rows
```

Average freight charges per order:

```
ala> var weekly_avg_fre = weekly_df.groupBy(col("WeekNumber")).agg(round(sum(col("order_freight_value"))/count("*"),3).as("Avg_Rev_per_Order")).orderBy("WeekNumber
aekly_avg_fre: org.apache.spark.sql.Dataset[org.apache.spark.sql.Row] = [WeekNumber: int, Avg_Rev_per_Order: double]
 ala> weekly avg fre.show()
.
WeekNumber|Avg_Rev_per_Order|
                                    21.869
                                   22.24
21.588
21.808
                                   21.738
21.306
21.871
22.863
21.608
            10
11
12
13
14
15
16
17
18
                                    21.881
21.646
                                    22.257
21.792
                                    21.503
21.064
                                    21.162
                                    21.233
             19
                                    21.614
             20|
                                    22.053
 ly showing top 20 rows
```

converting the three 'timestamp' columns to unified format:

```
var df2 - date_df.withColumn("Neekdumber", weekofyear(col("PurchaseDate"))).withColumn("order_purchase_timestamp", to_timestamp(col("order_purchase_timestamp"), "dd/MM/yy H:mm")).withColumn("order_aproved_at", to_timestamp(col("order_delivered_customer_date"), "dd/MM/yy H:mm"))
% apache.spark.sql.DataFrame = [id: int, order_status: string ... 17 more fields]
id|order_status|order_products_value|order_freight_value|order_items_qty| customer_city|customer_state|customer_zip_code_prefix|product_name_length|product_description_length|product_photos_qty|review_score|orestamp| order_aproved_at|order_delivered_customer_date|PurchaseDate|ApproveOate|DeliveryDate|MeckNumber|
  1| delivered|

:56:00|2017-10-02 11:07:00|

2| delivered|

:41:00|2018-07-26 03:24:00|
                                                                                                                                          1| Luziania|
2017-10-02| 2017-10-02| 2017-10-10|
1| Joinville|
2018-07-24| 2018-07-26| 2018-08-07|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     201
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   50 l
                                                                                                               119.9
                                                                                                                                        | 27.16|

2018-08-07 15:27:00|

41.69|

2018-08-17 18:06:00|

| 17.92|

2017-12-02 00:28:00|

| 23.11|

2018-02-16 18:17:00|

| 23.02|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   892 l
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    511
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              2018-07-24
                                                                                                                                                                                                                       2018-07-24 | 2018-07-25 | 2018-08-07 | 2018-08-07 | 2018-08-07 | 2018-08-07 | 2018-08-07 | 2018-08-07 | 2018-08-17 | 1 | RIO DE JAMETRO] | 2017-11-18 | 2017-11-20 | 2018-02-13 | 2018-02-13 | 2018-02-13 | 2018-02-13 | 2017-07-09 | 2017-07-09 | 2017-07-09 | 2017-07-09 | 2017-07-09 | 2017-07-09 | 2017-07-09 | 2017-07-09 | 2017-07-09 | 2017-08-10 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-13 | 2017-08-1
                                                                                                                                                                                                                                                                                                                                                                                                 SC|
30|
ES|
32|
RJ|
46|
SP|
7|
SP|
48|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   1156
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              2018-08-08
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     40
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                2018-02-13
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   1591
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     692 l
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              2017-07-09
                                                                                                                                          23.02

2017-07-26 10:57:00

| 16.21

2017-04-13 13:25:00

| 19.82

| 2017.05-26 12:55:00
                                                                                                                   65.0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   370
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     893
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                2017-04-11
            delivered|
:00|2017-04-13 13:25:00|
         2:00]2017-04-13 13:25:00]
| delivered|
9:00]2017-05-16 13:22:00]
| delivered|
9:00]2017-01-25 02:50:00]
| delivered|
5:00]2017-07-29 12:05:00]
                                                                                                                                          | 19.82|
| 2017-05-26 12:55:00|
| 51.14|
| 2017-02-02 14:08:00|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  2017-01-23
                                                                                                                   2017-02-02 14:08:00|
56.99| 16.13|
2017-08-16 17:14:00|
599.0| 15.69|
2017-05-29 11:18:00|
62.0| 2017-07-19 14:04:00|
250.0| 35.02|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       641
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  2017-07-29
            delivered|
:00|2017-05-16 19:50:00|
                                                                                                                  599.0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     125
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    1686
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 2017-05-16
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      882|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    50|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 2017-07-13
              delivered|
00|2017-07-13 20:10:00|
delivered|
```

Average time taken to approve the orders (order approved – order purchased):

Average order delivery time:

c. Total freight charges:

```
var weeklyFreight_df = weekly_df.groupBy("WeekNumber").agg(round(sum(col("order_freight_value")),2).as("Total_Freight_Charges")).orderBy("WeekNumber").show()
|
|WeekNumber|Total Freight Charges|
                              31469.52
                             41699.72
42398.0
                             42656.18
45150.2
45765.64
           4|
5|
6|
7|
8|
9|
                              45054.13
                             49041.98
52119.31
          10
                              49145.77
                             48075.61
52992.88
44608.61
         11 |
12 |
13 |
14 |
15 |
16 |
                             47091.58
                             42233.19
                              49080.8
          17
                              49920.09
                              54567.61
          19
                              60562.31
                              61682.45
          20
only showing top 20 rows
 eeklyFreight_df: Unit = ()
```

d. Distribution of freight charges in each city: