

camila-emr-workspace

March 16, 2024

```
[1]: from pyspark.sql import SparkSession
      from pyspark.sql.functions import col
```

VBox()

Starting Spark application

<IPython.core.display.HTML object>

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

SparkSession available as 'spark'.

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

```
[3]: spark = SparkSession.builder.appName("CamilaRedfinDataAnalysis").getOrCreate()
```

VBox()

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

```
[4]: # extrai do bucket emr-camila-raw-data para o início das análises
      redfin_data = spark.read.csv("s3://emr-camila-raw-data/city_market_tracker.
      ↳tsv000.gz", header=True, inferSchema=True, sep= "\t")
```

VBox()

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

```
[5]: redfin_data.show(3)
```

VBox()

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

```
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
```



```

-----+-----+-----+-----+
-----+-----+-----+-----+
---+-----+-----+-----+-----+
-----+-----+-----+

```

only showing top 3 rows

```
[6]: #visualiza o schema
redfin_data.printSchema()
```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',
  ↳ layout=Layout(height='25px', width='50%'),...
```

root

```

|-- period_begin: date (nullable = true)
|-- period_end: date (nullable = true)
|-- period_duration: integer (nullable = true)
|-- region_type: string (nullable = true)
|-- region_type_id: integer (nullable = true)
|-- table_id: integer (nullable = true)
|-- is_seasonally_adjusted: string (nullable = true)
|-- region: string (nullable = true)
|-- city: string (nullable = true)
|-- state: string (nullable = true)
|-- state_code: string (nullable = true)
|-- property_type: string (nullable = true)
|-- property_type_id: integer (nullable = true)
|-- median_sale_price: double (nullable = true)
|-- median_sale_price_mom: double (nullable = true)
|-- median_sale_price_yoy: double (nullable = true)
|-- median_list_price: double (nullable = true)
|-- median_list_price_mom: double (nullable = true)
|-- median_list_price_yoy: double (nullable = true)
|-- median_ppsf: double (nullable = true)
|-- median_ppsf_mom: double (nullable = true)
|-- median_ppsf_yoy: double (nullable = true)
|-- median_list_ppsf: double (nullable = true)
|-- median_list_ppsf_mom: double (nullable = true)
|-- median_list_ppsf_yoy: double (nullable = true)
|-- homes_sold: integer (nullable = true)
|-- homes_sold_mom: double (nullable = true)
|-- homes_sold_yoy: double (nullable = true)
|-- pending_sales: integer (nullable = true)
|-- pending_sales_mom: double (nullable = true)
|-- pending_sales_yoy: double (nullable = true)
|-- new_listings: integer (nullable = true)
|-- new_listings_mom: double (nullable = true)
|-- new_listings_yoy: double (nullable = true)

```

```

|-- inventory: integer (nullable = true)
|-- inventory_mom: double (nullable = true)
|-- inventory_yoy: double (nullable = true)
|-- months_of_supply: double (nullable = true)
|-- months_of_supply_mom: double (nullable = true)
|-- months_of_supply_yoy: double (nullable = true)
|-- median_dom: integer (nullable = true)
|-- median_dom_mom: integer (nullable = true)
|-- median_dom_yoy: integer (nullable = true)
|-- avg_sale_to_list: double (nullable = true)
|-- avg_sale_to_list_mom: double (nullable = true)
|-- avg_sale_to_list_yoy: double (nullable = true)
|-- sold_above_list: double (nullable = true)
|-- sold_above_list_mom: double (nullable = true)
|-- sold_above_list_yoy: double (nullable = true)
|-- price_drops: double (nullable = true)
|-- price_drops_mom: double (nullable = true)
|-- price_drops_yoy: double (nullable = true)
|-- off_market_in_two_weeks: double (nullable = true)
|-- off_market_in_two_weeks_mom: double (nullable = true)
|-- off_market_in_two_weeks_yoy: double (nullable = true)
|-- parent_metro_region: string (nullable = true)
|-- parent_metro_region_metro_code: integer (nullable = true)
|-- last_updated: timestamp (nullable = true)

```

```

[7]: #print os nomes das colunas
      redfin_data.columns

```

VBox()

```

FloatProgress(value=0.0, bar_style='info', description='Progress:',
  ↳layout=Layout(height='25px', width='50%'),...

```

```

['period_begin', 'period_end', 'period_duration', 'region_type',
'region_type_id', 'table_id', 'is_seasonally_adjusted', 'region', 'city',
'state', 'state_code', 'property_type', 'property_type_id', 'median_sale_price',
'median_sale_price_mom', 'median_sale_price_yoy', 'median_list_price',
'median_list_price_mom', 'median_list_price_yoy', 'median_ppsf',
'median_ppsf_mom', 'median_ppsf_yoy', 'median_list_ppsf',
'median_list_ppsf_mom', 'median_list_ppsf_yoy', 'homes_sold', 'homes_sold_mom',
'homes_sold_yoy', 'pending_sales', 'pending_sales_mom', 'pending_sales_yoy',
'new_listings', 'new_listings_mom', 'new_listings_yoy', 'inventory',
'inventory_mom', 'inventory_yoy', 'months_of_supply', 'months_of_supply_mom',
'months_of_supply_yoy', 'median_dom', 'median_dom_mom', 'median_dom_yoy',
'avg_sale_to_list', 'avg_sale_to_list_mom', 'avg_sale_to_list_yoy',
'sold_above_list', 'sold_above_list_mom', 'sold_above_list_yoy', 'price_drops',
'price_drops_mom', 'price_drops_yoy', 'off_market_in_two_weeks',
'off_market_in_two_weeks_mom', 'off_market_in_two_weeks_yoy',
'parent_metro_region', 'parent_metro_region_metro_code', 'last_updated']

```

```
[8]: df_redfin = redfin_data.select(['period_end', 'period_duration', 'city', 'state', 'property_type', 'median_sale_price', 'median_ppsf', 'homes_sold', 'inventory', 'months_of_supply', 'median_dom', 'sold_above_list', 'last_updated'])
df_redfin.show(3)
```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',
layout=Layout(height='25px', width='50%'),...
```

```
+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+
|period_end|period_duration|city|state|
property_type|median_sale_price|
median_ppsf|homes_sold|inventory|months_of_supply|median_dom|
sold_above_list|last_updated|
+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+
|2021-02-28|30|Fair Lawn|New Jersey|Single Family
Res...|542500.0|278.1372476167421|20|69|
3.5|91|0.0|2024-03-10 14:36:40|
|2019-04-30|30|Elyria|Ohio|Multi-Family
(2-4...|30000.0|18.610421836228287|3|7|
2.3|130|0.0|2024-03-10 14:36:40|
|2020-08-31|30|Northwest Harwich|Massachusetts|Single Family
Res...|625000.0|340.65934065934067|11|14|
1.3|26|0.18181818181818182|2024-03-10 14:36:40|
+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+
only showing top 3 rows
```

```
[9]: #verifica o número total de linhas
print(f"Total number of rows: {df_redfin.count()}")
```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',
layout=Layout(height='25px', width='50%'),...
```

Total number of rows: 5245047

```
[10]: from pyspark.sql.functions import isnull
#Conta valores nulos em cada coluna
#usamos uma list comprehension
```

```

null_counts = [df_redfin.where(isnull(col_name)).count() for col_name in
↳df_redfin.columns]
null_counts

```

VBox()

```

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

```

[0, 0, 0, 0, 0, 6066, 69125, 5646, 413699, 334643, 69227, 36370, 0]

```

[11]: #Exibe os resultados
for i, col_name in enumerate(df_redfin.columns):
    print(f"{col_name}: {null_counts[i]} null values")

```

VBox()

```

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

```

```

period_end: 0 null values
period_duration: 0 null values
city: 0 null values
state: 0 null values
property_type: 0 null values
median_sale_price: 6066 null values
median_ppsf: 69125 null values
homes_sold: 5646 null values
inventory: 413699 null values
months_of_supply: 334643 null values
median_dom: 69227 null values
sold_above_list: 36370 null values
last_updated: 0 null values

```

```

[12]: # Verifica se há valores faltantes em todo o DataFrame
remaining_count = df_redfin.na.drop().count()

print(f"Number of missing rows: {df_redfin.count() - remaining_count}")

```

VBox()

```

FloatProgress(value=0.0, bar_style='info', description='Progress:',
↳layout=Layout(height='25px', width='50%'),...

```

Number of missing rows: 501546

```

[13]: print(f"Total number of remaining rows: {remaining_count}")

```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳layout=Layout(height='25px', width='50%'),...
```

Total number of remaining rows: 4743501

```
[14]: #remove NA e conta o número total de linhas restantes  
df_redfin = df_redfin.na.drop()  
print(f"Total number of rows: {df_redfin.count()}")
```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳layout=Layout(height='25px', width='50%'),...
```

Total number of rows: 4743501

```
[15]: # Conta os valores nulos em cada coluna para confirmar se removemos todos os  
↳na(tem q ficar tudo zero)  
null_counts = [df_redfin.where(isnull(col_name)).count() for col_name in  
↳df_redfin.columns]  
null_counts
```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳layout=Layout(height='25px', width='50%'),...
```

[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]

```
[16]: from pyspark.sql.functions import year, month  
  
#Extrai o ano de period_end e salva em uma nova coluna "period_end_yr"  
df_redfin = df_redfin.withColumn("period_end_yr", year(col("period_end")))  
  
#Extrai mês de period_end e salva em uma nova coluna "period_end_month"  
df_redfin = df_redfin.withColumn("period_end_month", month(col("period_end")))
```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳layout=Layout(height='25px', width='50%'),...
```

```
[17]: # Elimine as colunas period_end e last_updated  
df_redfin = df_redfin.drop("period_end", "last_updated")
```

VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳layout=Layout(height='25px', width='50%'),...
```

```
[18]: df_redfin.show(3)
```


VBox()

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳ layout=Layout(height='25px', width='50%'),...
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+  
-----+-----+-----+-----+-----+-----+-----+-----+  
-----+-----+-----+-----+  
|period_duration|          city|          state|  
property_type|median_sale_price|  
median_ppsf|homes_sold|inventory|months_of_supply|median_dom|  
sold_above_list|period_end_yr|period_end_month|  
+-----+-----+-----+-----+-----+-----+-----+-----+  
-----+-----+-----+-----+-----+-----+-----+-----+  
-----+-----+-----+-----+  
|          30|          Fair Lawn|          New Jersey|Single Family Res...|  
542500.0| 278.1372476167421|          20|          69|          3.5|          91|  
0.0|          2021|          2|  
|          30|          Elyria|          Ohio|Multi-Family (2-4...|  
30000.0|18.610421836228287|          3|          7|          2.3|          130|  
0.0|          2019|          4|  
|          30|Northwest Harwich|Massachusetts|Single Family Res...|  
625000.0|340.65934065934067|          11|          14|          1.3|  
26|0.18181818181818182|          2020|          8|  
+-----+-----+-----+-----+-----+-----+-----+-----+  
-----+-----+-----+-----+-----+-----+-----+-----+  
-----+-----+-----+-----+  
only showing top 3 rows
```

```
[19]: from pyspark.sql.functions import when  
  
#altera o número do mês para o respectivo nome do mês.  
  
df_redfin = df_redfin.withColumn("period_end_month",  
                                when(col("period_end_month") == 1, "January")  
                                .when(col("period_end_month") == 2, "February")  
                                .when(col("period_end_month") == 3, "March")  
                                .when(col("period_end_month") == 4, "April")  
                                .when(col("period_end_month") == 5, "May")  
                                .when(col("period_end_month") == 6, "June")  
                                .when(col("period_end_month") == 7, "July")  
                                .when(col("period_end_month") == 8, "August")  
                                .when(col("period_end_month") == 9, "September")  
                                .when(col("period_end_month") == 10, "October")  
                                .when(col("period_end_month") == 11, "November")  
                                .when(col("period_end_month") == 12, "December")  
                                .otherwise("Unknown")  
                                )
```

```
VBox()
```

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳ layout=Layout(height='25px', width='50%'),...
```

```
[20]: df_redfin.show(3)
```

```
VBox()
```

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳ layout=Layout(height='25px', width='50%'),...
```

```
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+
|period_duration|          city|          state|
property_type|median_sale_price|
median_ppsf|homes_sold|inventory|months_of_supply|median_dom|
sold_above_list|period_end_yr|period_end_month|
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+
|          30|          Fair Lawn|          New Jersey|Single Family Res...|
542500.0| 278.1372476167421|          20|          69|          3.5|          91|
0.0|          2021|          February|
|          30|          Elyria|          Ohio|Multi-Family (2-4...|
30000.0|18.610421836228287|          3|          7|          2.3|          130|
0.0|          2019|          April|
|          30|Northwest Harwich|Massachusetts|Single Family Res...|
625000.0|340.65934065934067|          11|          14|          1.3|
26|0.18181818181818182|          2020|          August|
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+
only showing top 3 rows
```

```
[22]: #salva dataframe final transformado no outro bucket s3 como um arquivo parquet.
s3_bucket = "s3://emr-camila-transformed-data/redfin_data.parquet"
df_redfin.write.mode("overwrite").parquet(s3_bucket)
```

```
VBox()
```

```
FloatProgress(value=0.0, bar_style='info', description='Progress:',  
↳ layout=Layout(height='25px', width='50%'),...
```