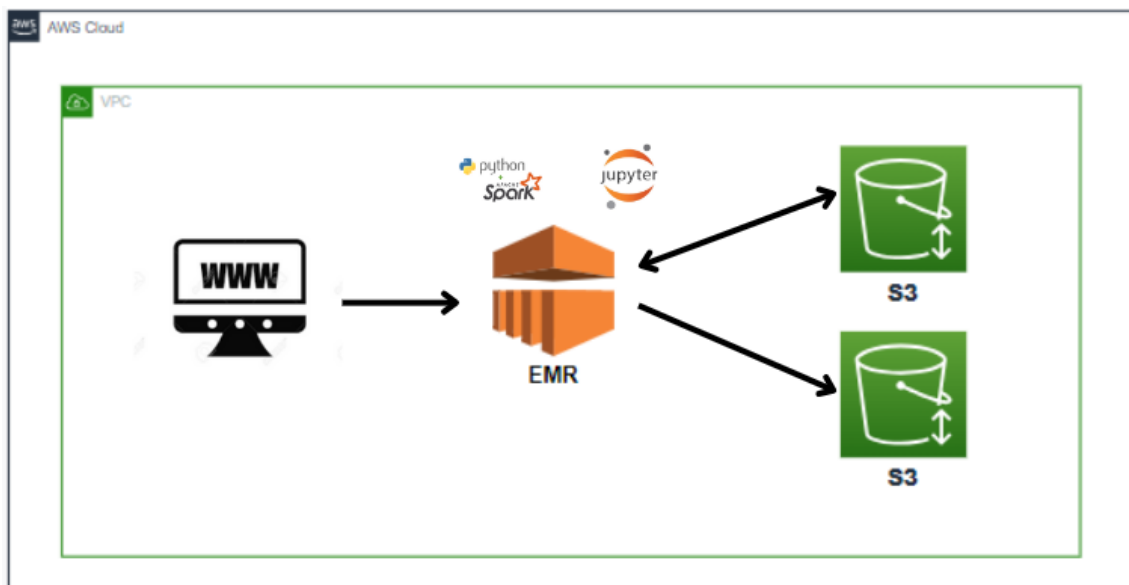# Creation of ETL Pipeline in EMR clusters using PySpark, Jupyter Notebook and S3 Buckets

## 1.Objective

In this project I built a data pipeline in the AWS cloud, from scratch, including environment configurations, creation of users, networks, subnets, EC2, security groups, EMR clusters, creation of S3 buckets, transformation, and cleaning of raw data, using PySpark, until the database containing ready-to-use data is made available to the end user. The data source is the website of Redfin, one of the largest real estate agencies in the US and Canada, which freely makes real estate market data available for metropolitan areas, cities, neighborhoods and zip codes. Below are the step-by-step instructions for each step.
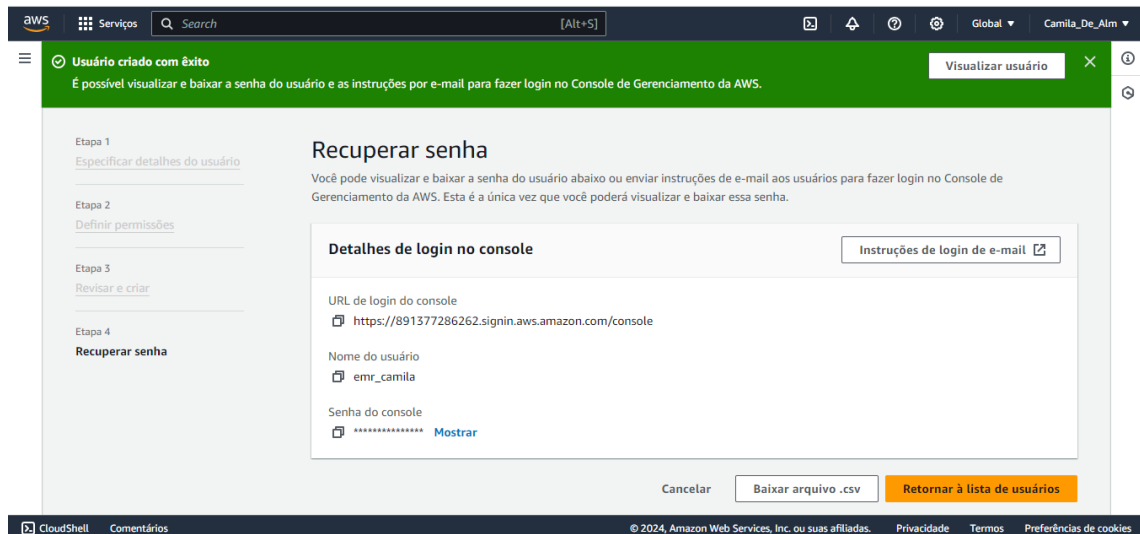
Demonstration of the final architecture flowchart on AWS:



## 2.Creation of the Administrator user

You will need to have an active AWS – Amazon World Services account to perform the following items. Due to security reasons, the root user, who is the user responsible for creating the AWS account, must create another administrator user so that in the event of an invasion of the account, it remains secure. For this project, the administrator user was created in AWS IAM (Identity and Access Management) which is an AWS service that helps you control access to AWS resources securely. Logged in as the root user, I created the credentials and keypairs for the
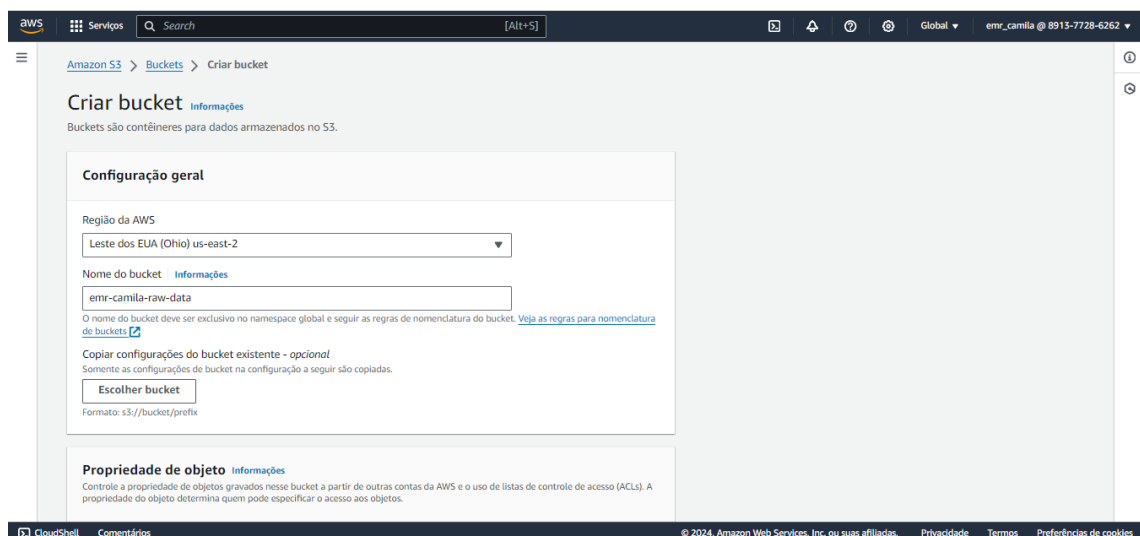
administrator user, and enabled the necessary permissions to use the EMR and S3 Buckets.



For security reasons, I also downloaded the credentials in csv, so I could recover them in the future, in case of lost or forgotten password. After following the security practices recommended by AWS, I logged out of the Root user and logged in with the new Administrator user.
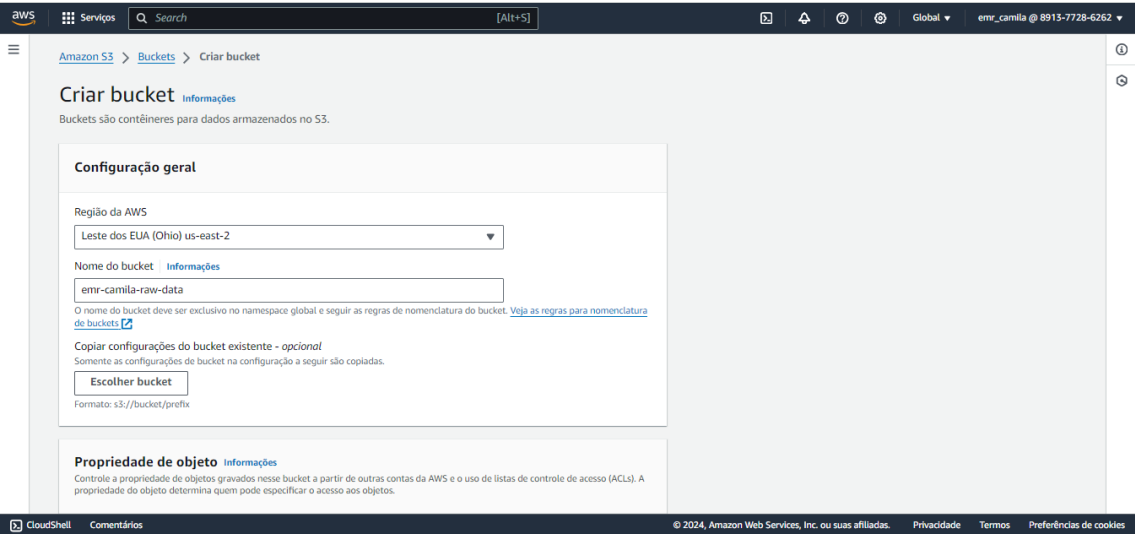
## 3. Creating Buckets in AWS S3

Using the Administrator user, I created two S3 buckets, the first to receive the data extraction from the source located on the Redfin website and the second bucket to receive the transformed and cleaned data. In the AWS console, search for S3, then create the bucket. The bucket name must be globally unique and follow the naming rules determined by AWS, as shown in the image below.



The region chosen was us-east of the USA, due to the cost being the most recommended, but due to latency, it is recommended that the Region chosen to be
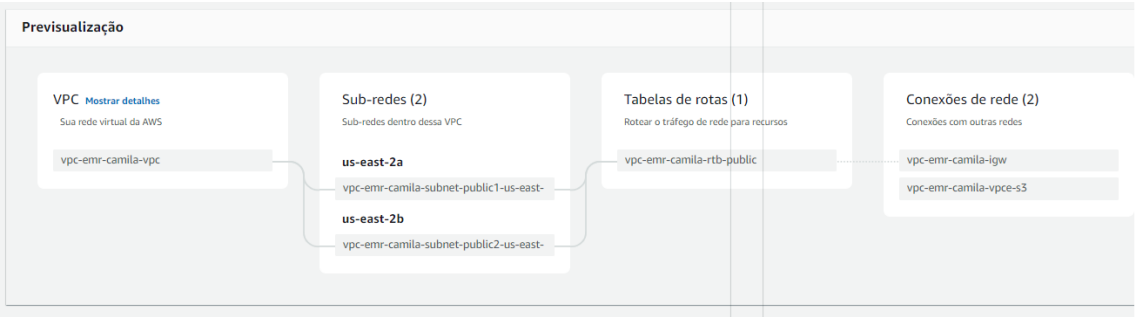
the same as where your user is located. Additionally, it is possible to determine properties, tags, and versioning, but I kept the default options suggested by AWS. For the second bucket, I followed the same procedure as above.



## 4.Creation of VPC and Subnets

VPC is a Virtual Private Cloud, that is, a virtual private cloud, it creates an isolated virtual private network in AWS, protecting my EMR cluster and my data against unauthorized access and encrypts data at rest and in transit between Amazon S3 and Amazon EMR to ensure information confidentiality. In the console, search for VPC, then click "Create VPC". I chose a name for the VPC, we must create the VPC, subnets, route tables and determine the availability zones, gateway and endpoint . I configured the availability zones as two, because if there is a disaster and one region becomes unavailable, data transit will be redirected to the second, I set the same value for public networks and none for private networks, because it will be internal data. The endpoint will be the S3 Gateway.

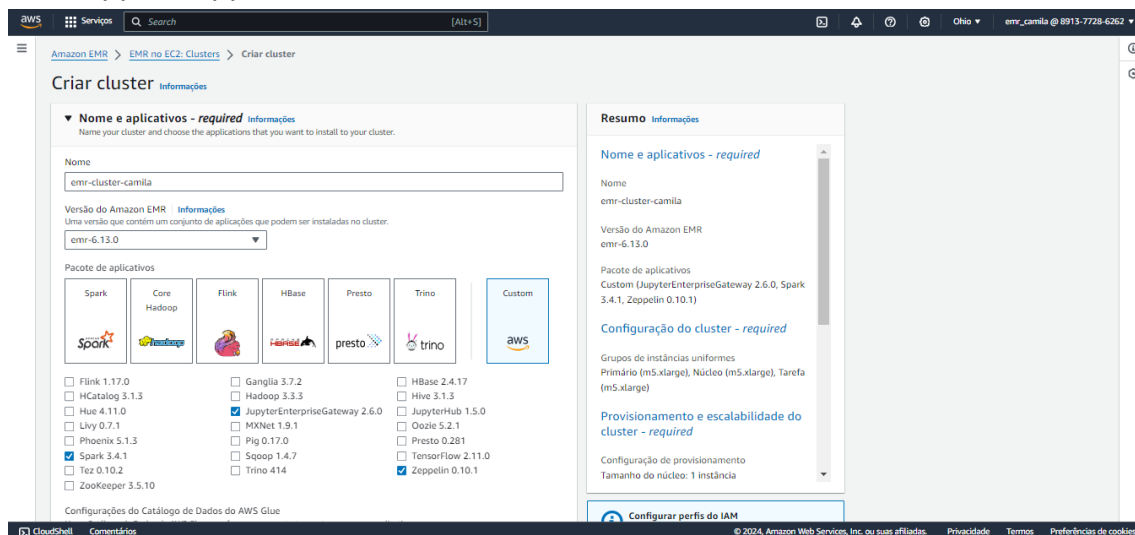Next, visualization of the entire VPC flow.

## 5. EMR Clusters

After configuring the user, environment, and security, I created the AWS EMR (Elastic Map Reduce) cluster, which is an AWS service that facilitates the execution of big data frameworks, such as Apache Hadoop and Apache Spark, to process and analyze large volumes of data, in this project I used Apache Spark. In AWS EMR clusters, nodes are EC2 instances that come together to provide the processing and storage power needed to run your big data jobs. There are three main types of nodes in an EMR cluster:

- **Node or Primary Node:** It is responsible for coordinating all work in the cluster, managing tasks, monitoring the health of the cluster, and tracking the status of jobs.
- **Core Nodes:** They perform data processing tasks and there can be multiple core nodes to distribute and execute jobs in parallel.
- **Task Nodes (Optional):** Provide additional processing power to perform CPU-intensive tasks.

In this project, I created the Primary Node and the Secondary Node, due to the frequency and size of the load. The EMR already comes with main frameworks, so it can be configured directly in the console, which reduces the cost. The setup chosen was version 6.13.0, I didn't choose the most recent one, I selected the Spark, Jupyter and Zeppelin applications, as shown below.



The group of instances chosen was uniform, to use AWS EC2 on demand or spot, due to economies of scale, the type chosen was M5.xlarge, which contains 4v core and 16Gib of memory, for the core I chose the same configurations and I removed the option to have a task cluster , because it is optional. Amazon EC2 M5 instances are powered by the cloud's fastest Intel Xeon Scalable processors with an all -core

turbo frequency of up to 4.5 GHz. Additionally, M5zn instances feature 100 Gbps networking capacity, making them ideal for computing-intensive applications with high consumption of network resources, but the instance can be changed depending on the size of the load, for performance measures I chose the M5.Xlarge.



For the "Cluster provisioning and scalability" section, regarding scalability I chose the option managed by EMR, because this way the EMR adapts according to the variation in the data flow. The minimum cluster size (number of instances) was three, because one is for the Main Node and the second and third are for the Central Node. When demand is high, the maximum number of the cluster will be 10 EC2 instances, while for the Central Node it will be 8. Provisioning configuration is set to two, which would be the initial capacity of the core instance.

In the "Networks" section, I selected the VPC that I had previously created, searched, selected the VPC and AWS automatically filled in the information about the VPC and Subnets.



In the topic "Cluster termination and node replacement", I selected "Automatically terminate the cluster after downtime", because it is a way to reduce cost, so it will be terminated in times of downtime. I selected the option "Use termination protection" which refers to EC2, thus deleting the clusters, but keeping the stored data.

In the section, "Cluster logs", I allowed AWS to create a bucket to store the EMR cluster logs.



Regarding "EC2 security and key pair configuration", the recommended practice is to create a keypair to access via SSH, however I created it directly through the console. Click on the create key pair button, opened another tab , gave it a name, chose the RSA pair and the format. pem , I saved the generated file.

Then select the keypair in the security configuration as below.



In the "IAM Profiles" section, I selected the option choose a service profile, that's why I created the VPC previously, I selected the default security group each time we create a VPC a security group is created automatically.

▼ **Perfis do Identity and Access Management (IAM) - *required*** Informações
Escolha ou crie um perfil de serviço e um perfil de instância para as instâncias do EC2 no cluster.

**Perfil de serviço do Amazon EMR** Informações
O perfil de serviço é um perfil do IAM que o Amazon EMR assume para provisionar recursos e executar ações de nível de serviço com outros serviços da AWS.

○ Escolha um perfil de serviço existente
Selecione um perfil de serviço padrão ou um perfil personalizado com políticas do IAM anexadas para que o cluster possa interagir com outros serviços da AWS.

● Escolha um perfil de serviço
Deixe que o Amazon EMR crie um novo perfil de serviço para que você possa conceder e restringir o acesso a recursos em outros serviços da AWS.

**Recursos de rede**
Já adicionamos os recursos que você configurou na seção Redes. Escolha a VPC, a sub-rede e os grupos de segurança que o perfil de serviço pode acessar.

Nuvem privada virtual (VPC)

*Escolha uma ou mais VPCs* ▼

vpc-emr-camila-vpc ✕
vpc-0a78b076a8755ae72

Sub-rede

*Escolha uma ou mais sub-redes* ▼

vpc-emr-camila-subnet-public1-us-east-2a ✕
subnet-0d794e8e8058361f9

Grupo de segurança

*Escolha um ou mais grupos de segurança* ▼

default ✕
sg-0a3ac6005ed529794

After in the "Instance Profile" section, I chose an instance profile, gave access to all buckets for reading and writing, which is why the buckets were created before the clusters.

## Perfil de instância do EC2 para o Amazon EMR

O perfil de instância atribui um perfil a cada instância do EC2 em um cluster. O perfil de instância deve especificar um perfil que possa acessar os recursos para as etapas e ações de bootstrap.

○ **Escolha um perfil de instância existente**
Selecione um perfil padrão ou um perfil de instância personalizado com políticas do IAM anexadas para que o cluster possa interagir com seus recursos no Amazon S3.

● **Escolha um perfil de instância**
Deixe que o Amazon EMR crie um novo perfil de instância para que você possa especificar um conjunto personalizado de recursos para acesso no Amazon S3.

**Acesso ao bucket do S3** | Informações

○ Buckets ou prefixos específicos do S3 em sua conta  Informações
Escolha os buckets ou prefixos que você deseja que esse perfil de instância acesse.

● Todos os buckets do S3 nessa conta com acesso de leitura e gravação
Conceda ao perfil de instância acesso a todos os buckets que tiverem acesso de leitura e gravação habilitado em sua conta.

## Função personalizada de ajuste de escala automático - *opcional*

Quando uma regra personalizada de ajuste de escala automático é acionada, o Amazon EMR assume essa função para adicionar e encerrar instâncias do EC2. Saiba mais ↗

Função personalizada de ajuste de escala automático

| Escolher perfil do IAM ▼ | ⟳ | Criar perfil do IAM ↗ |

Afterwards, I clicked on "Create cluster", which can take up to 7 minutes to complete.



## 6. Configuring Jupyter Notebook to use PySpark

Now it is possible to configure Jupyter Notebook to use PySpark, for this we have to configure a Studio. On the left, in the EMR Studio section, click on basic concepts, image below. After creating studio, a bucket was created to be able to attach Jupyter to the EMR, so the codes written in the notebook will be executed in the EMR and in the S3 bucket via programming. To do this, I needed to create a custom bucket and assign it to my user, first I gave it a name, then I created a bucket just for my studio and associated it with my administrator user.

Afterwards, in the "Networks and security" section, I selected the VPC and Subnet that I had already created, remembering that the administrator user must have full access permission to AWS S3 and be an Administrator User, as mentioned in step two, after clicking "Create Studio". Now we can create the Jupyter notebook, click on the URL of the created Studio, as shown below.



I clicked on "Create workspace" chose a name, the other settings remain the same for network and storage.

VPC that you had already created.



S3 bucket created specifically for Jupyter.

I clicked on Studio below, and it opened Jupyter Notebook in a new tab.



Now it is necessary to attach Jupyter to the cluster that we created, and we must select the EC2 option.

I refreshed the page and the message below appeared under "EMR Computing".





Afterwards, I went back to my workspace, double-clicked on the notebook and chose the kernel as Pyspark .



## 7. Code Explanation

In the first cell, we imported Spark, in the second, I loaded directly from the Redfin page and stored it in the first bucket created, emr -camila- raw -data, the bucket that was empty, already received the load, as shown below.

After processing and cleaning, I saved the transformed data in the second bucket I created, as shown below, we have the file and the success log. The file was saved in Parquet because of the size, but it could have been in csv, but aiming for better storage and lower cost, I chose the Parquet extension.



## 8. Complete Jupyter Notebook

Next, I attached the Jupyter Notebook with all the codes, transformations and changes, the cells have commented code. It is important to remember that the database provided by RedFin was mostly clean and organized, however, column adjustments, checking NA, Null values and changing variable names were applied. All the basis was taken from the AWS website and the idea for the project was taken from the Telespectra channel.

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%'),…
```

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

```
+-----------------+---------------+-----------------+---------------+-------------+---------------+-------------+--------+-------------+------------+---------------+---------------+-----------------+-------------------+-----------------+---------------+-----------------+-------------------+------------------+----------+---------------+-------------+-------------+---------------+-------------+-------------+-------------+-------------+------------------+------------------+-----------------+----------+-------------+-------------+--------------+------------------+------------------+-------------+------------------+----------------+-----------+-----------------+-----------+------------------+-------------------+-------------------------+-------------------------+----------------------+-------------+
|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|
+-----------------+---------------+-----------------+---------------+-------------+---------------+-------------+--------+-------------+------------+---------------+---------------+-----------------+-------------------+-----------------+---------------+-----------------+-------------------+------------------+----------+---------------+-------------+-------------+---------------+-------------+-------------+-------------+-------------+------------------+------------------+-----------------+----------+-------------+-------------+--------------+------------------+------------------+-------------+------------------+----------------+-----------+-----------------+-----------+------------------+-------------------+-------------------------+-------------------------+----------------------+-------------+
|     2021-02-01|2021-02-28|             30|      place|            6|   5779|f|        Fair Lawn, NJ|     Fair Lawn|  New Jersey|        NJ|Single Family Res…|               6|         542500.0|  0.21378230227094752|0.35964912280701755|         469000.0|  -0.14711765775595564|0.04454342984409809| 278.1372476167421|-0.03125719269388…|-0.00459122296016…|291.19741100323625|-0.02531486920188…|-0.02420798775354…|        20|-0.4117647058823529|0.4285714285714286|           14|-0.26315789473684215|0.2727272727272727|
```

```
25|0.7857142857142858|-0.3902439024390244|          69|
0.0|-0.34905660377358494|                3.5|              1.5|
-4.1|          91|          8|          -2|
1.0|-0.0121431598112657|0.020401919432001314|               0.0|-0.20588235294
117646|-0.14285714285714285|0.08695652173913043|-0.07246376811594202|
0.03978671041837571|   0.21428571428571427|        0.10902255639097744|
-0.05844155844155843|      New York, NY|
35614|2024-03-10 14:36:40|
|   2019-04-01|2019-04-30|           30|     place|          6|    6504|
f|         Elyria, OH|        Elyria|          Ohio|      OH|Multi-Family
(2-4…|          4|          30000.0|  -0.7014925373134329|
-0.6557165399512266|         110000.0|  -0.04264577893820709|
-0.06542056074766356|18.610421836228287|  -0.6026777101693752|
-0.5572373301855638|
81.52173913043478|0.004286982260352179|0.005327375837772763|        3|
0.5|             0.5|             2|             null|          1.0|
4|          3.0|0.33333333333333326|          7|          0.0|
-0.5|           2.3|  -1.2000000000000002|              -4.7|      130|
-149|          33|
0.860358532597008|-0.02600566975054…|0.007997142263547308|           0.0|
0.0|             0.0|             null|             null|
null|             1.0|             null|
1.0|      Cleveland, OH|                17460|2024-03-10 14:36:40|
|   2020-08-01|2020-08-31|           30|     place|          6|    24459|
f|Northwest Harwich…|Northwest Harwich|Massachusetts|       MA|Single Family
Res…|          6|          625000.0|   0.3789299503585217|
-0.14617486338797814|         649000.0|   0.26019417475728157|
0.42668718399648276|340.65934065934067|   0.12348060103162162|
0.31284153005464477| 393.0817610062893|   0.2546073635603787|
0.18042570872759556|         11|0.2222222222222232|2.6666666666666665|
8|  -0.4285714285714286|             1.0|
14|0.5555555555555556|0.07692307692307687|         14|0.07692307692307687|
-0.7407407407407407|              1.3|-0.09999999999999987|              -16.7|
26|          -51|          12|0.9771452130091354|0.022599268740304268|-0.0270
9829717127…|0.18181818181818182|  0.07070707070707072|  -0.4848484848484848|
0.2857142857142857|-0.09890109890109894|0.007936507936507908|
0.5|        0.4285714285714286|                0.5|Barnstable Town, MA|
12700|2024-03-10 14:36:40|
+-----------+----------+--------------+----------+--------------+--------+---
----------------+-----------------+----------------+-------------+------
---+-----------------+-------------+----------------+----------------+------
-+----------------+---------------+----------------+--------------+--------------
-----+-------------+-----------------+----------------+-------------+------------
----+-------------+----------------+--------+----------------+------------+-
--------------+--------------+---------------+----------------+-------
----+-------------+--------------+---------------+----------------+-----
-------------+--------------+----------------+----------------+------
----+-------------+------------+----------------+----------------+-----
```

```
--------------+-----------------+------------------+------------------+--
----------------+-----------------+-----------------+------------------
---+------------------------+------------------------+-----------------+
---------------------------+-----------------+
```
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]:
```python
#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]

```
[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]:
```python
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%'),…
```