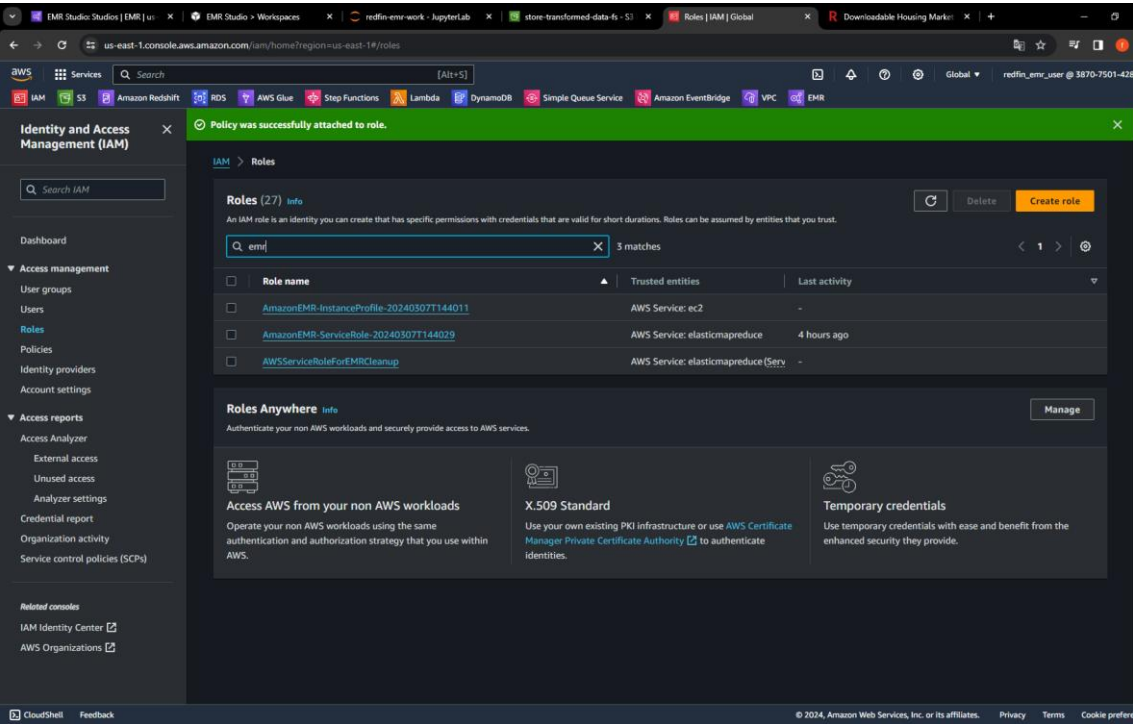
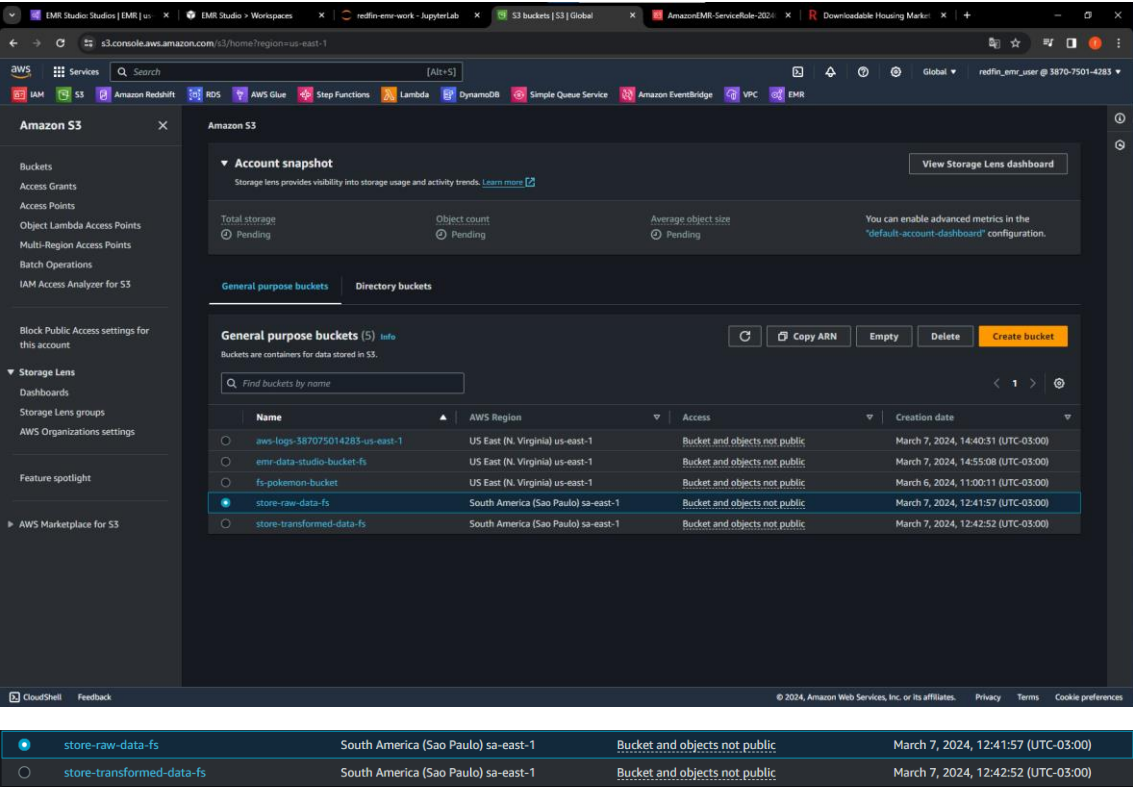


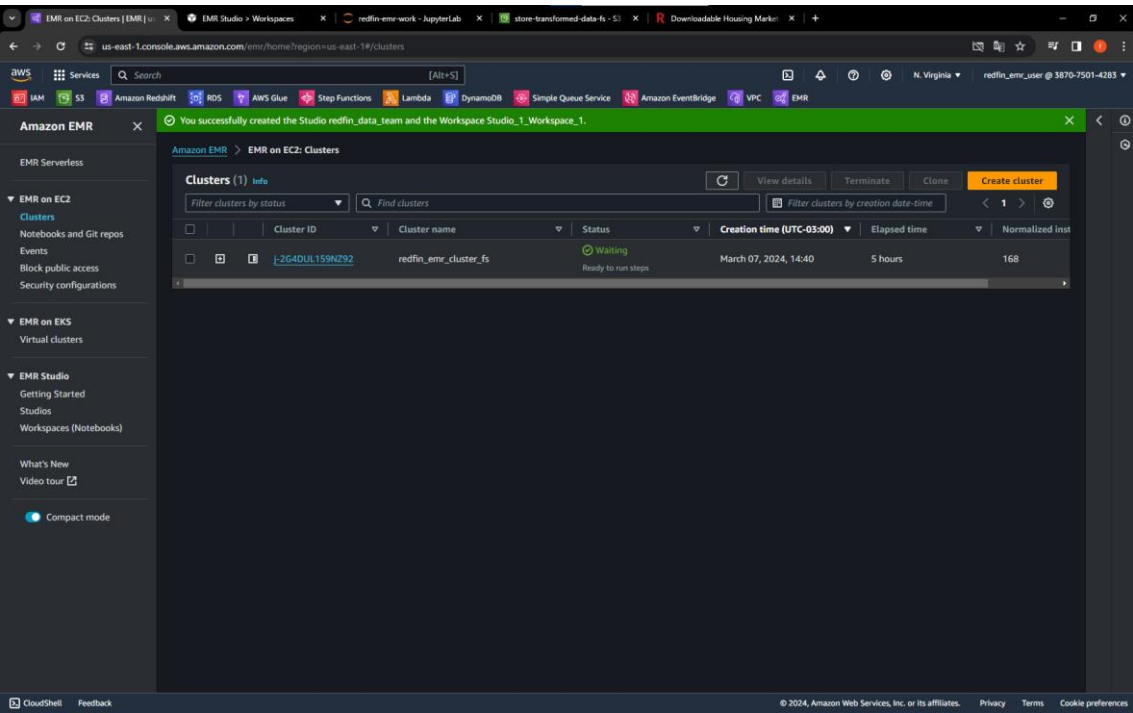
IAM Roles and permissions:



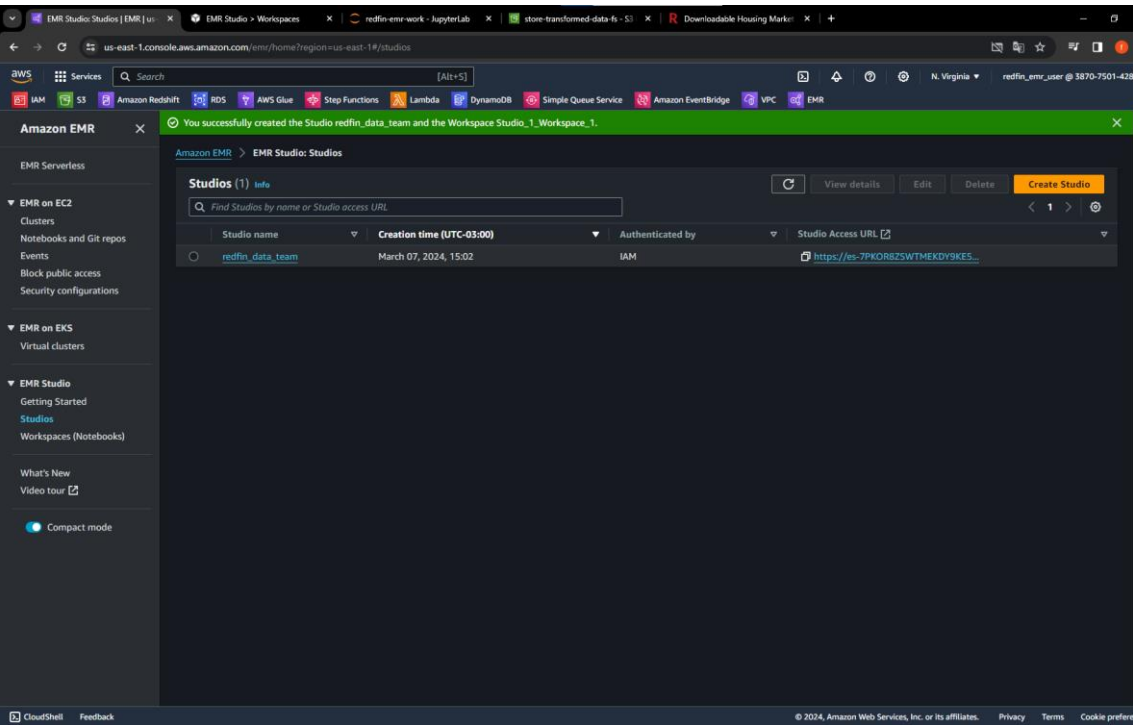
S3 Buckets:



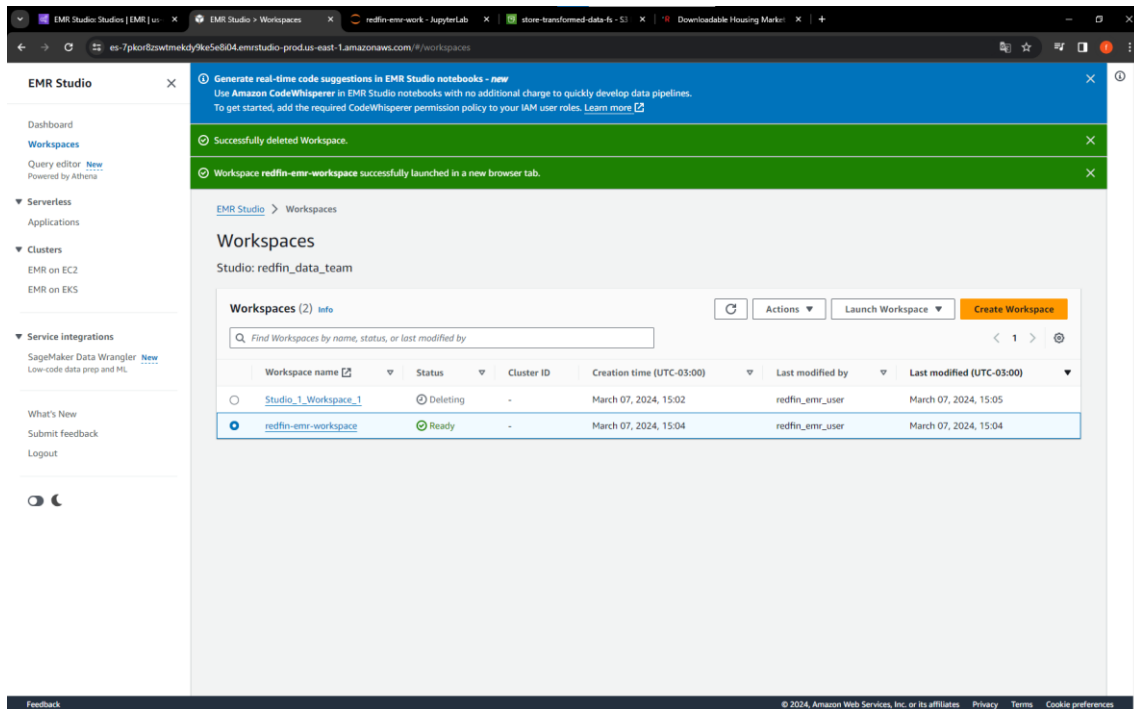
EMR on EC2 Cluster:



EMR Studio:

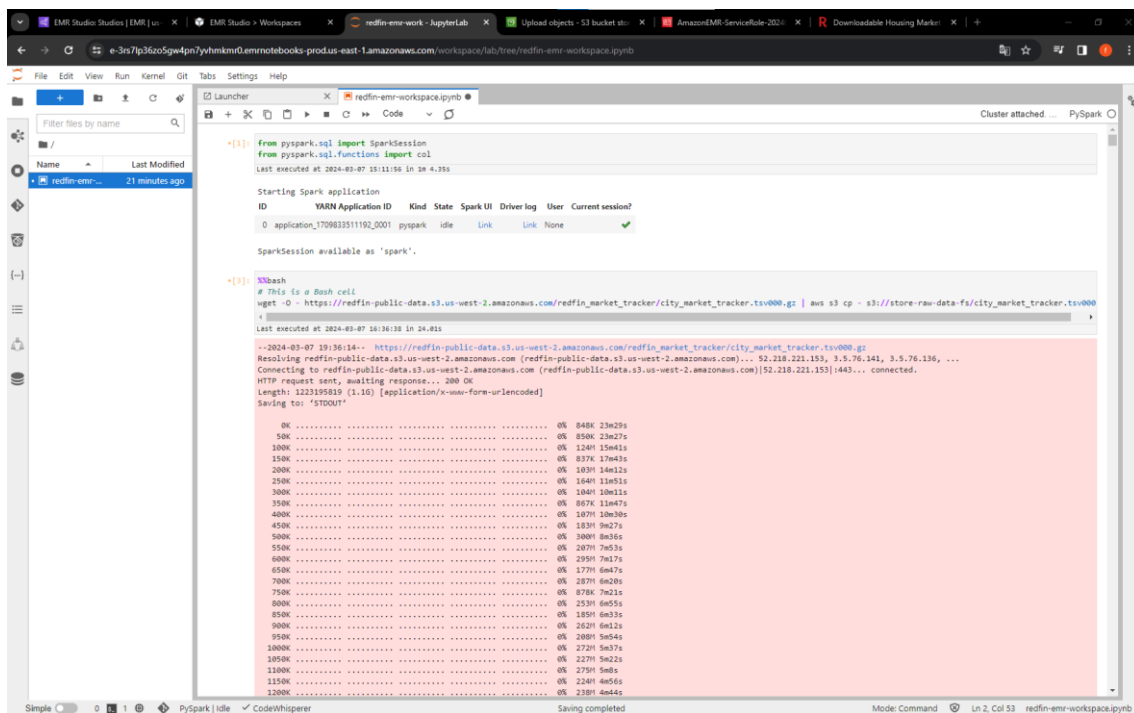


EMR Workspace:

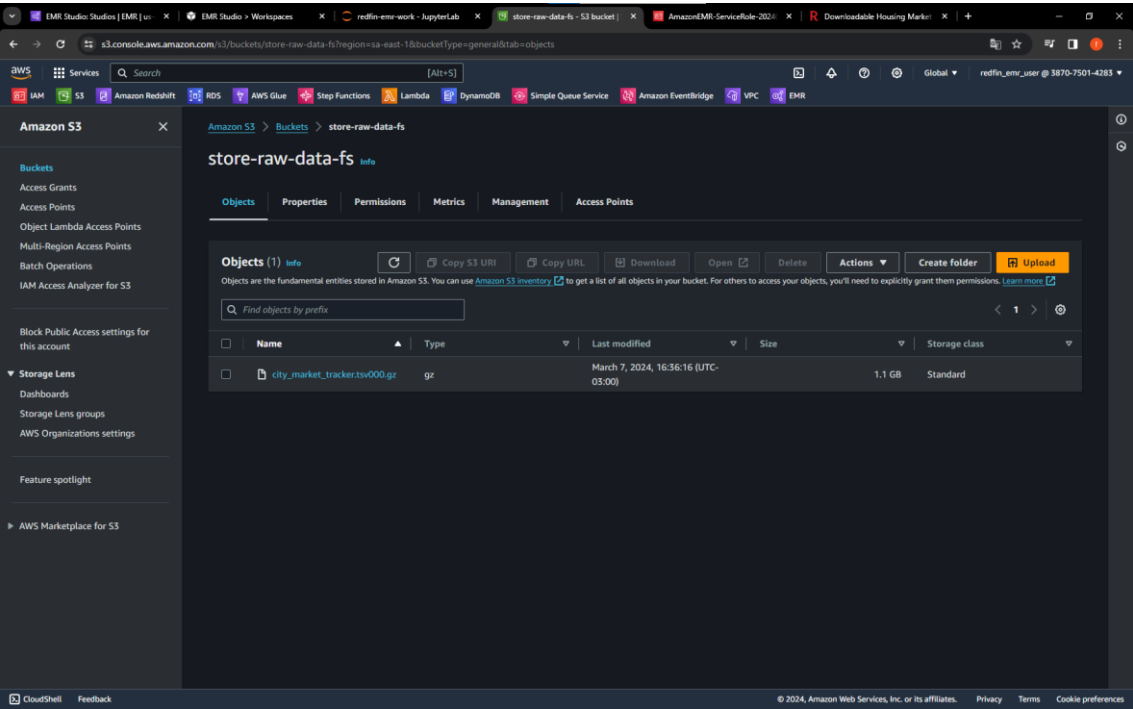


PySpark Notebook:

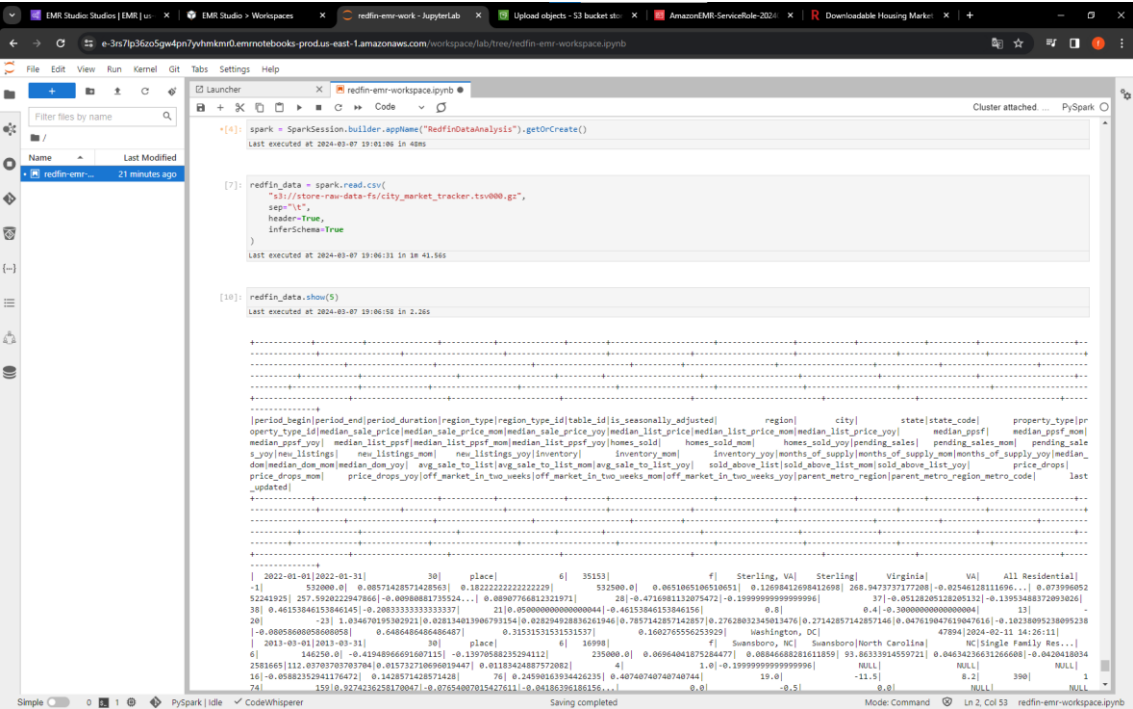
(Api request)



S3 Raw data bucket, file uploaded:



Pyspark Notebook transformations:



EMR Studio | [MR] | us | x | EMR Studio | Workspaces | x | redefin-emr-workspace | JupyterLab | x | Upload objects - S3 bucket sto... | x | AmazonEMR-ServiceRole-2024... | x | Downloadable Housing Marke... | x | +

e-3rs7p36zo5gw4pn7yhnkm0.emrnotebooks-prod-us-east-1.amazonaws.com/workspace/lab/tree/redefin-emr-workspace.pyrb

File Edit View Run Kernel Git Tabs Settings Help

Launcher redefin-emr-workspace.pyrb Cluster attached... PySpark

Filter files by name

Name Last Modified

redefin-emr... 22 minutes ago

redefin\_data.printSchema()

Last executed at 2024-03-07 13:07:03 in 348s

```
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)
```

Simple 0 1 PySpark | Idle CodeWhisperer Saving completed Mode: Command Ln 2, Col 53 redefin-emr-workspace.pyrb

EMR Studio | [MR] | us | x | EMR Studio | Workspaces | x | redefin-emr-workspace | JupyterLab | x | Upload objects - S3 bucket sto... | x | AmazonEMR-ServiceRole-2024... | x | Downloadable Housing Marke... | x | +

e-3rs7p36zo5gw4pn7yhnkm0.emrnotebooks-prod-us-east-1.amazonaws.com/workspace/lab/tree/redefin-emr-workspace.pyrb

File Edit View Run Kernel Git Tabs Settings Help

Launcher redefin-emr-workspace.pyrb Cluster attached... PySpark

Filter files by name

Name Last Modified

redefin-emr... 22 minutes ago

redefin\_data.columns

Last executed at 2024-03-07 13:07:04 in 33ms

```
[ '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_yoy', 'off_market_in_two_weeks', 'off_market_in_two_weeks_mom', 'off_market_in_two_weeks_yoy', 'parent_metro_region_metro_code', 'last_updated' ]
```

[15]: selected\_columns = [ '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' ]

Last executed at 2024-03-07 13:08:19 in 33ms

df\_redefin = redefin\_data.select(selected\_columns)

Last executed at 2024-03-07 13:08:09 in 239ms

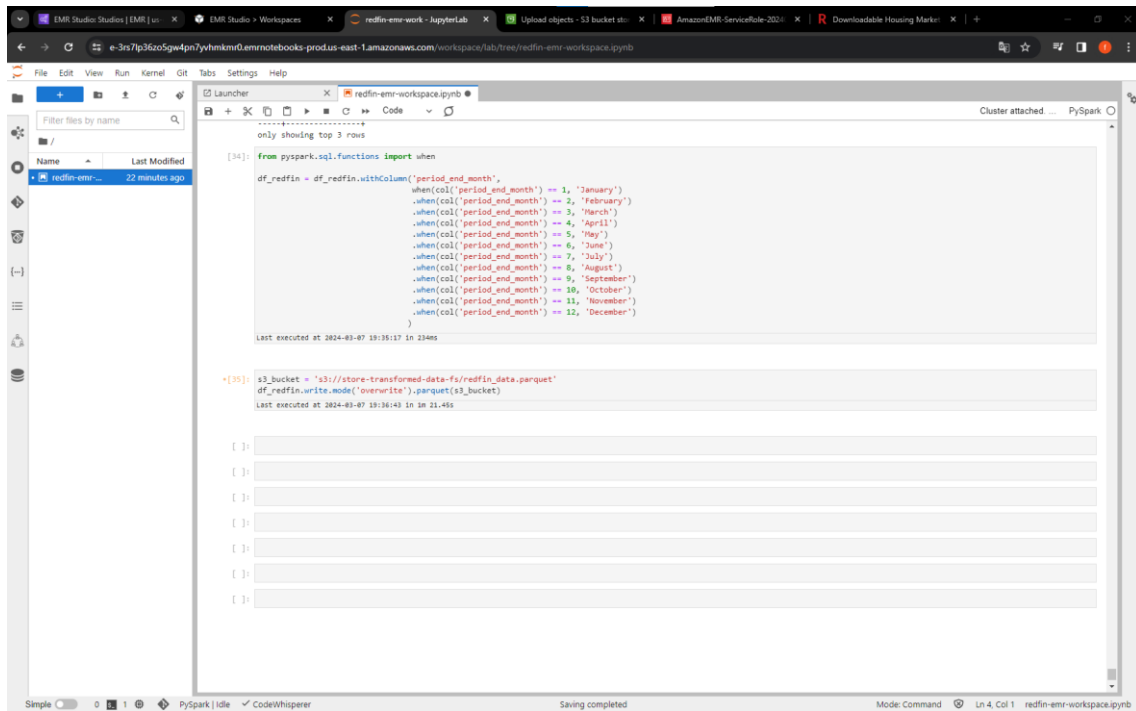
df\_redefin.show(5)

Last executed at 2024-03-07 13:09:17 in 9.28s

[period_end period_duration	city	state	property_type	median_sale_price	median_ppsf	homes_sold	inventory	months_of_supply	median_dom	sold_above_list
2022-01-31	30	Sterling	Virginia	All Residential	532000.0	268.9473737177208	28	21	0.8	13
2024-02-11 14:26:11	30	Suwanboro	North Carolina	Single Family Res...	146250.0	93.86333914559721	41	76	19.0	390

Simple 0 1 PySpark | Idle CodeWhisperer Saving completed Mode: Command Ln 2, Col 53 redefin-emr-workspace.pyrb





## S3 Transformed bucket:

