

python™

```

redfin_analytics.py X
C: > Users > DELL > Desktop > Un_Uploaded Projects > Redfin_Snowpipe > redfin_analytics.py > transform_data
11 # s3 buckets
12 target_bucket_name = 'redfin-transformed-data'
13
14 url_by_city = 'https://redfin-public-data.s3.us-west-2.amazonaws.com/redfin_market_tracker/city_market_tracker.tsv000.gz'
15
16 def extract_data(**kwargs):
17     url = kwargs['url']
18     df = pd.read_csv(url, compression='gzip', sep='\t')
19     now = datetime.now()
20     date_now_string = now.strftime("%d%m%Y%H%M%S")
21     file_str = 'redfin_data_' + date_now_string
22     df.to_csv(f"{file_str}.csv", index=False)
23     output_file_path = f"/home/ubuntu/{file_str}.csv"
24     output_list = [output_file_path, file_str]
25     return output_list
26
27 def transform_data(task_instance):
28     data = task_instance.xcom_pull(task_ids="tsk_extract_redfin_data")[0]
29     object_key = task_instance.xcom_pull(task_ids="tsk_extract_redfin_data")[1]
30     df = pd.read_csv(data)
31
32     # Remove commas from the 'city' column
33     df['city'] = df['city'].str.replace(',', '')
34     cols = ['period_begin', 'period_end', 'period_duration', 'region_type', 'region_type_id', 'table_id',
35            'is_seasonally_adjusted', 'city', 'state', 'state_code', 'property_type', 'property_type_id',
36            'median_sale_price', 'median_list_price', 'median_ppsf', 'median_list_ppsf', 'homes_sold',
37            'inventory', 'months_of_supply', 'median_dom', 'avg_sale_to_list', 'sold_above_list', 'parent_metro_region_metro_code', 'last_updated']
38     df = df[cols]
39     df = df.dropna()
40

```



Apache
Airflow

Press **shift** + **/** for Shortcuts

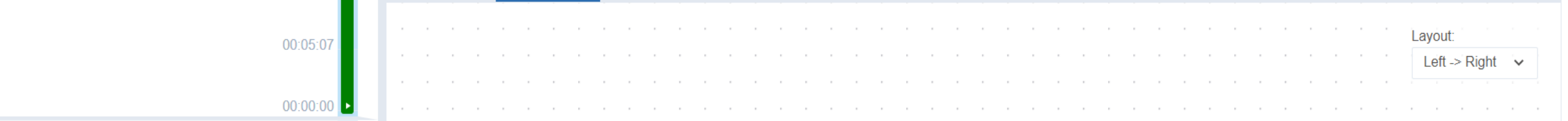
deferred failed queued removed restarting running scheduled shutdown skipped success up_for_reschedule up_for_retry upstream_failed no_status



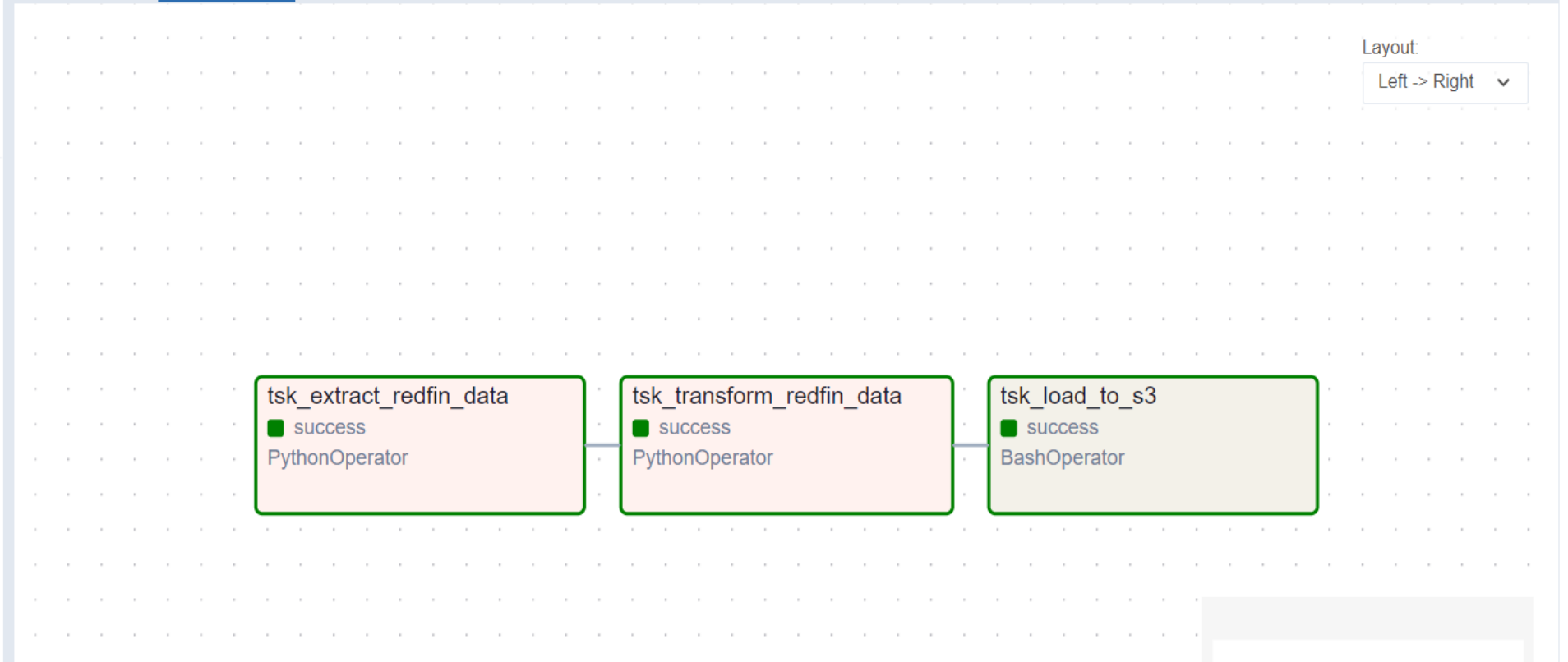
DAG **redfin_analytics_dag** / Run **2024-09-22, 00:00:00 UTC**

Clear Mark state as...

Details Graph Gantt Code Event Log



tsk_extract_redfin_data	■
tsk_transform_redfin_data	■
tsk_load_to_s3	■





Amazon
S3

Instances | EC2 | us-west-2

redfin-rawdata - S3 bucket

EC2 Instance Connect

redfin_analytics_dag - Gr...

RealEstateScript - Snowfl...

us-west-2.console.aws.amazon.com/s3/buckets/redfin-rawdata?region=us-west-2&bucketType=general&tab=objects

☆Z⋮

aws

Services

Search

[Alt+S]

Oregon

Zyad_Ahmed

Amazon S3

Buckets

redfin-rawdata

redfin-rawdata

Info

Objects

Properties

Permissions

Metrics

Management

Access Points

Objects (1)

Info

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

<1>

	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	<div><div></div><div>redfin_data_04102024130246.csv</div></div>	csv	October 4, 2024, 16:10:53 (UTC+03:00)	2.9 GB	Standard

CloudShell

Feedback

© 2024, Amazon Web Services, Inc. or its affiliates.

Privacy

Terms

Cookie preferences

Instances | EC2 | us-west-2

redfin-transformed-data

EC2 Instance Connect

redfin_analytics_dag - Gri

RealEstateScript - Snowfl

us-west-2.console.aws.amazon.com/s3/buckets/redfin-transformed-data?region=us-west-2&bucketType=general&tab=objects

☆ Z

aws

Services

Search [Alt+S]

Oregon

Zyad_Ahmed

Amazon S3

Buckets

redfin-transformed-data

redfin-transformed-data

Info

Objects

Properties

Permissions

Metrics

Management

Access Points

Objects (1) Info

Copy S3 URI

Copy URL

Download

Open

Delete

Actions

Create folder

Upload

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

Find objects by prefix

< 1 >

	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	<div><div></div><div>redfin_data_04102024130246.csv</div></div>	csv	October 4, 2024, 16:10:40 (UTC+03:00)	953.9 MB	Standard

CloudShell

Feedback

© 2024, Amazon Web Services, Inc. or its affiliates.

Privacy

Terms

Cookie preferences



snowflake®

EC2 Instance Connect

Instances | EC2 | us-west-

EC2 Instance Connect

redfin_analytics_dag - Gri

RealEstateScript - Snowfl

app.snowflake.com/foyvqcu/ohb03577/w4BoRICrwhXq#query

RealEstateScript

Databases

Worksheets

Search objects

REDFIN_DATABASE_1

- EXTERNAL_STAGE_SCHEMA
- FILE_FORMAT_SCHEMA
- INFORMATION_SCHEMA
- PUBLIC
- REDFIN_SCHEMA
- SNOWPIPE_SCHEMA

SNOWFLAKE

SNOWFLAKE_SAMPLE_DATA

ACCOUNTADMIN

COMPUTE_WH (X-Small)

Share

REDFIN_DATABASE_1.SNOWPIPE_SCHEMA

Settings

Code Versions

69

70

71

72

73

74

75

76

77

SELECT *

FROM redfin_database_1.redfin_schema.redfin_table LIMIT 5;

SELECT COUNT(*) FROM redfin_database_1.redfin_schema.redfin_table

-- DESC TABLE redfin_database.redfin_schema.redfin_table;

Results

Chart

	PERIOD_BEGIN	PERIOD_END	PERIOD_DURATION	REGION_TYPE	REGION_TYPE_ID	TABLE_ID	IS_SEASONALLY_ADJUSTE
1	2017-09-01	2017-09-30	30	place	6	29470	f
2	2023-03-01	2023-03-31	30	place	6	38334	f
3	2020-07-01	2020-07-31	30	place	6	37598	f
4	2022-09-01	2022-09-30	30	place	6	14794	f

Ask Copilot

EC2 Instance Connect

Instances | EC2 | us-west-

EC2 Instance Connect

redfin_analytics_dag - Gri

RealEstateScript - Snowfl

app.snowflake.com/foyvqcu/ohb03577/w4BoRICrwhXq#query

☆ ⬇️ Z ⋮

RealEstateScript

+

▼

Databases

Worksheets

Search objects

↻

▼

REDFIN_DATABASE_1

>

EXTERNAL_STAGE_SCHEMA

>

FILE_FORMAT_SCHEMA

>

INFORMATION_SCHEMA

>

PUBLIC

>

REDFIN_SCHEMA

>

SNOWPIPE_SCHEMA

>

SNOWFLAKE

>

SNOWFLAKE_SAMPLE_DATA

REDFIN_DATABASE_1.SNOWPIPE_SCHEMA

Settings

Code Versions

⋮

69

70

71

72

73

74

75

76

77

SELECT *

FROM redfin_database_1.redfin_schema.redfin_table LIMIT 5;

SELECT COUNT(*) FROM redfin_database_1.redfin_schema.redfin_table

-- DESC TABLE redfin_database.redfin_schema.redfin_table;

Results

Chart

	COUNT(*)
1	4331055

Query Details

...

Query duration

76ms

Rows

1

Query ID

01b7778a-0003-45e8-0...

Show more

Ask Copilot

ZA



Power BI

File Home Insert Modeling View Optimize Help

Share

Clipboard

Paste

Cut

Copy

Format painter

Data

Get data

Excel workbook

OneLake data hub

SQL Server

Enter data

Dataverse

Recent sources

Queries

Transform data

Refresh

Insert

New visual

Text box

More visuals

Calculations

New visual calculation

New measure

Quick measure

Sensitivity

Share

Publish

Copilot

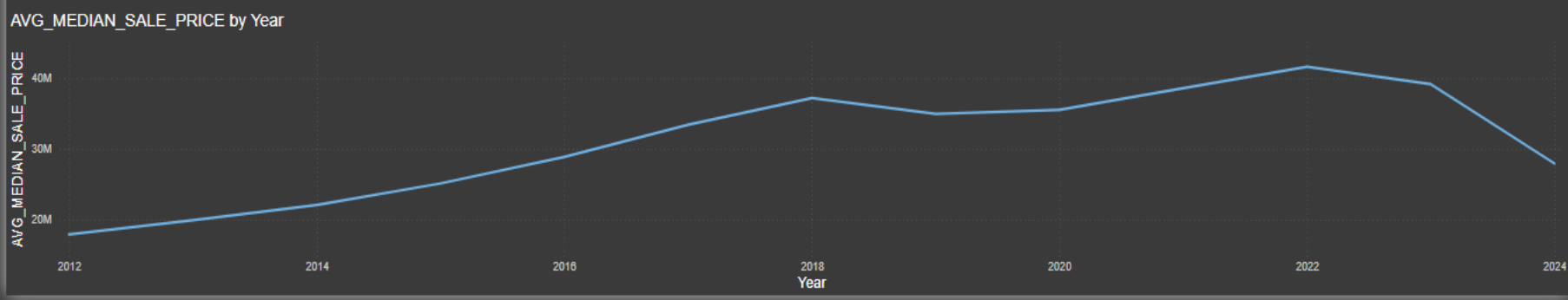
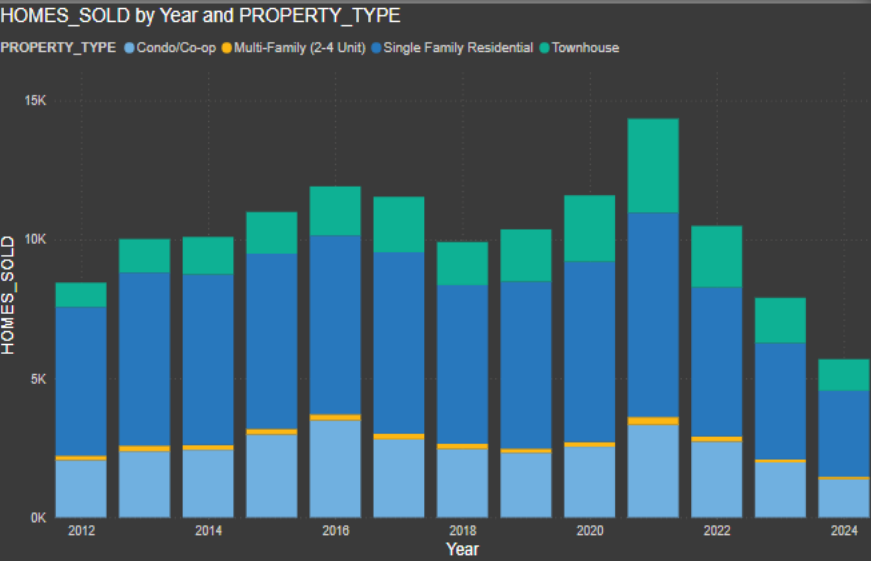
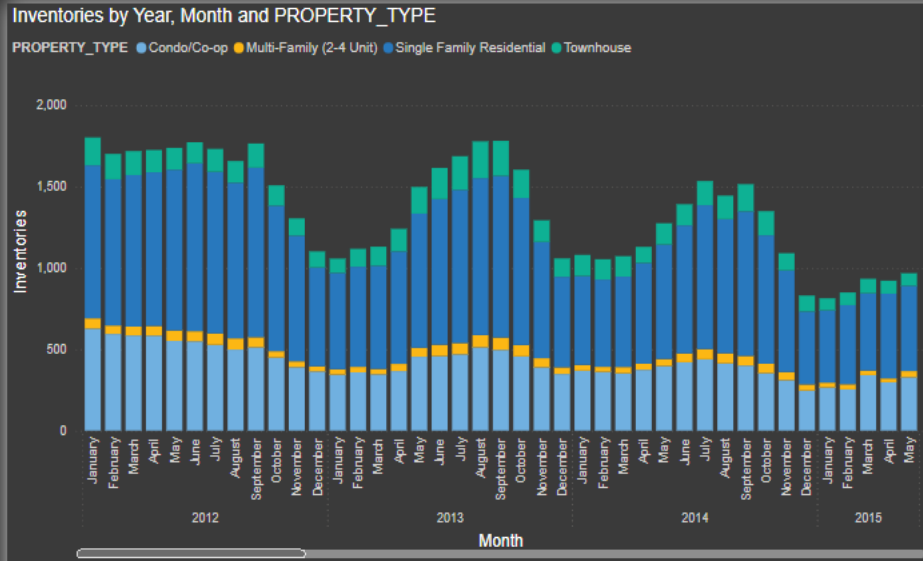
Visualizations

Grid

Table

Map

DAX



Visualizations

Data

Search

REDFIN_TABLE