# Objective:

Perform Analysis on the Real Estate Data and Retrieve Meaningful Information from it.

# Approach:

- Import the CSV file into Databricks.
- Analyze the data.
- Extract Information
- Create a Report

# Goal:

Extract meaningful data from the records and create an appropriate record.

## 1. Uploading to Databricks

Firstly, the file is uploaded to Databricks to be analyzed. On doing so, we can see the file in the dbfs along with its file path.
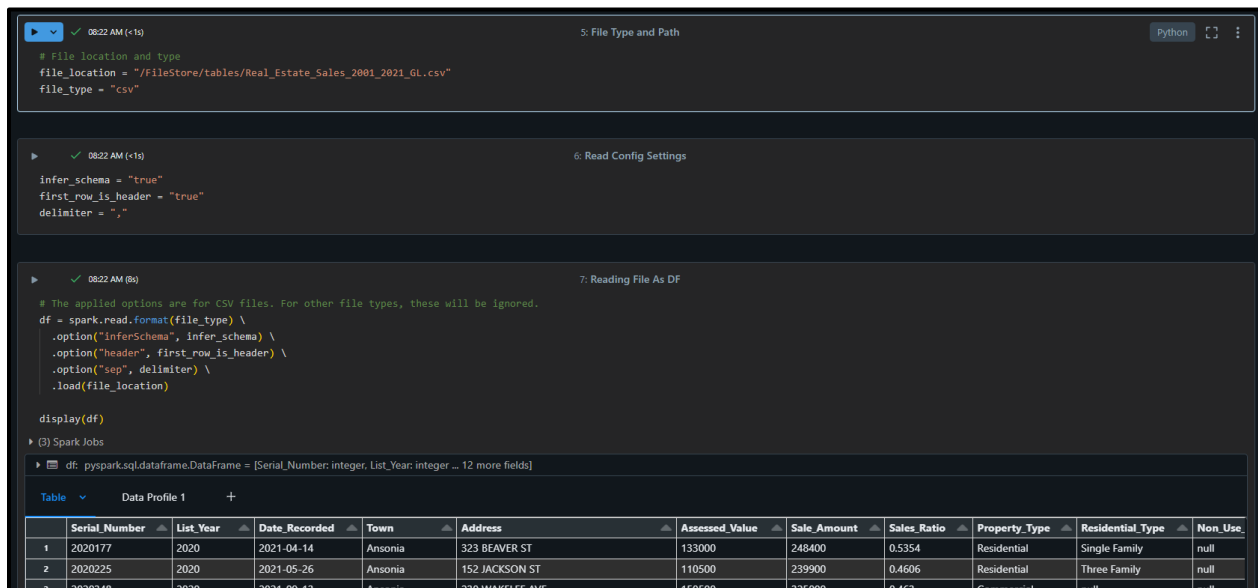


## 2. Reading the CSV

The file is then read using **spark.read** into a dataframe. Required configurations to read the file are also applied.

## 3. Saving DF to Table

The dataframe is then saved as a table in Databricks.



## 4. Analyze Table Content

Firstly, column data types and amount of records are read.



While checking number of records according to name, an anomaly is detected where the name of the town is ***Unknown***. The data is then further explored.

## 5. Town Name Anomaly

Since it is the only anomaly, the record is further explored. On doing so address can be retrieved. Googling the address resulted in the Town Name: 'East Hampton'. The record in updated accordingly.





Upon checking the town names, we can now see the absence of the anomaly is the result,

## 6. Checking for Unique Identifier

On taking a glance of the data, Serial Number can be assumed as the Unique Identifier for records. A check is performed for unique identifier to see for presence of a random Serial Number through different Towns and its occurrence count. We can see it occurring 81 times. We can now assume the Serial Number is not a unique identifier for the records.

```sql
%sql
SELECT count(SERIAL_NUMBER), Town
FROM RealEstate_CSV
WHERE SERIAL_NUMBER = '200023'
GROUP BY Town;
```

20: Checking for a Serial Number in Different Towns

(2) Spark Jobs

▶ ▤ _sqldf: pyspark.sql.dataframe.DataFrame = [count(SERIAL_NUMBER): long, Town: string]

| | count(SERIAL_NUMBER) | Town |
|---|---|---|
| 1 | 1 | Bethlehem |
| 2 | 1 | Litchfield |
| 3 | 1 | Sterling |
| 4 | 1 | Windsor Locks |
| 5 | 1 | Woodbridge |
| 6 | 1 | Wolcott |
| 7 | 1 | Cromwell |

81 rows  |  2.22 seconds runtime

```sql
%sql
SELECT count(1)
FROM RealEstate_CSV
WHERE SERIAL_NUMBER = '200023';
```

21: Count of Records containg Serial Number 200023

(2) Spark Jobs

▶ ▤ _sqldf: pyspark.sql.dataframe.DataFrame = [count(1): long]

| | count(1) |
|---|---|
| 1 | 81 |

Now, a check is performed to see if Serial Number could be Unique for a particular town. On doing so, it results in multiple counts of the same Serial Number in the same Town. The Serial Number not being unique to a particular Town is also determined. It is also seen through all the towns in the data.



Now, a check is performed to see if Serial Number is Unique to a Particular Town for a Particular List Year. We can see it is unique except for 1 case. Upon further exploring it, the record is the exact duplicate.



On trying to delete one of the records, it is not supported.

```
%sql
DELETE FROM RealEstate_CSV
WHERE SERIAL_NUMBER = 70086
  AND List_Year = 2007
  AND Town = 'East Hampton'
  AND ROWID NOT IN (
    SELECT MIN(ROWID)
    FROM RealEstate_CSV
    WHERE SERIAL_NUMBER = 70086
      AND List_Year = 2007
      AND Town = 'East Hampton'
  );
```

AnalysisException: Multi-column In predicates are not supported in the DELETE condition.



```
%sql
WITH RankedRows AS (
    SELECT *,
        ROW_NUMBER() OVER (PARTITION BY SERIAL_NUMBER, List_Year, Town ORDER BY (SELECT NULL)) AS RowNum
    FROM RealEstate_CSV
    WHERE SERIAL_NUMBER = 70086 AND List_Year = 2007 AND Town = 'East Hampton'
)
DELETE FROM RankedRows
WHERE RowNum > 1;
```

AnalysisException: [UNSUPPORTED_FEATURE.TABLE_OPERATION] The feature is not supported: Table  does not support DELETE. Please check the current catalog and namespace to make sure the qualified table name is expected, and also check the catalog implementation which is configured by "spark.sql.catalog".

So, for now, it is left as is.



Except for this case, serial number can be a unique identifier for a property for a particular year.

## 7. Creating Unique Identifier

A unique Identifier for each record can be created by concatenating Serial_Number, List_Year and Town.



## 8. Property Types

On checking the property types, we can see 12 property types including null with count 0. This is because it cannot be counted directly.

So only null property types are calculated which is more than 10000.



## 9. Replacing null Property Type

Null Property Types are replaced with string value 'N/A'.



Null Property Types are not available anymore and data with N/A type is same as previous count.

## 10.        Residential Property Types

On checking data for containing value in the Residential Type column, there was presence of the same data in property type and residential type.





Since, those contained records in the Residential Type column, we can assume that Property Type can be changed to Residential.

And, on re-checking the conditions, we can see that the property type is categorized properly.



## 11.        Average Sales Ratio

To get data for average sales ratio, first, we need to get data for the properties that have been sold. Then, the condition for average sales ratio across all towns can be calculated.



We can see that the average sales are 42.64% above the average estimates value.

## 12.        Sales to Assessment Ratio by Town

The same condition above can be applied to calculate the ratio by town.



## 13.        Average Sale-To-Assessment Ratio against Total Average

The records can be visualized using appropriate graphs. The graph below illustrates the average sale-to-assessment ratio to the overall sale-to-assessment ratio. This can also be used to indicate Towns with unusual increment in sale amount to assessed value.

## 14.        Average Sales Ratio

We can also get the average sales ratio.



It can also be organized by Towns. It visualizes the Towns with the most rapid increase in the property valuations.

## 15.        Average Sales Amount by Town

The average Sales Amount by Town can also be visualized to see where in average the highest sale amount of any property by town.



## 16.        Property Type Distribution

We can visualize the Property Type Distribution in the available data.



We can see that most of the property type is Residential, and the next largest data is not available. The rest of the properties cannot be visualized properly. So, property types except Residential and N/A are visualized separately.

Also, the count of Residential and N/A Property Type is also visualized.

Also, property types except Residential and N/A are counted by Town to be visualized.



The same is also done with Residential and N/A property types.

This is done so that the number of different properties across the towns can be visualized by the count and distribution across the records.

## 17.        Residential Property Type Distribution

Residential Property Type Distribution can be visualized. This shows the types and distribution across the records.



## 18.        Residential Property Type Count

Residential Property Type Count can also be visualized across all the towns.

## 19.      Sale Amount to List Year

We can also visualize the average sale mount of properties according to different list years. We can see a massive rise in 2006 with a massive fall in 2008. That is when the housing market crashed. Soon after there was a gradual rise to 2014. Then a dip occurred in 2015 and decrease in 2016. Again, a massive spike occurred since then up to 2020 which again fell during the covid years.



We can also see the sale amount distribution across the years quarterly. This can help indicate quarterly changes in sale amount.

We can see the presence of a null record, on investigating, we can see the presence of an incomplete data.



It can be fixed simply by deleting the record.

## 20.        Chain Records

During data exploration, records indicating to a serial number is indicated. This indicated towards chained records. This can be used further for determining things like changes to assessment ratio when changes are made to the property itself. For now, the records are simply visualized on a table.

## 21.        Dashboards

From the above 3 dashboards can be created.

- AVG Dashboard

  https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/6332182203914549/2743257718242605/1006725448119626/latest.html


- Property Dashboard

  https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/6332182203914549/2743257718242605/1006725448119626/latest.html


- Sale/Year Dashboard

  https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/6332182203914549/2743257718242605/1006725448119626/latest.html


## 22.        Notebook

The notebook is the following:

https://databricks-prod-cloudfront.cloud.databricks.com/public/4027ec902e239c93eaaa8714f173bcfc/6332182203914549/2743257718242605/1006725448119626/latest.html