

## Deliverable 2

### GitHub Link:

<https://github.com/saikrishnasanda/bigdata>

### Data Understanding:

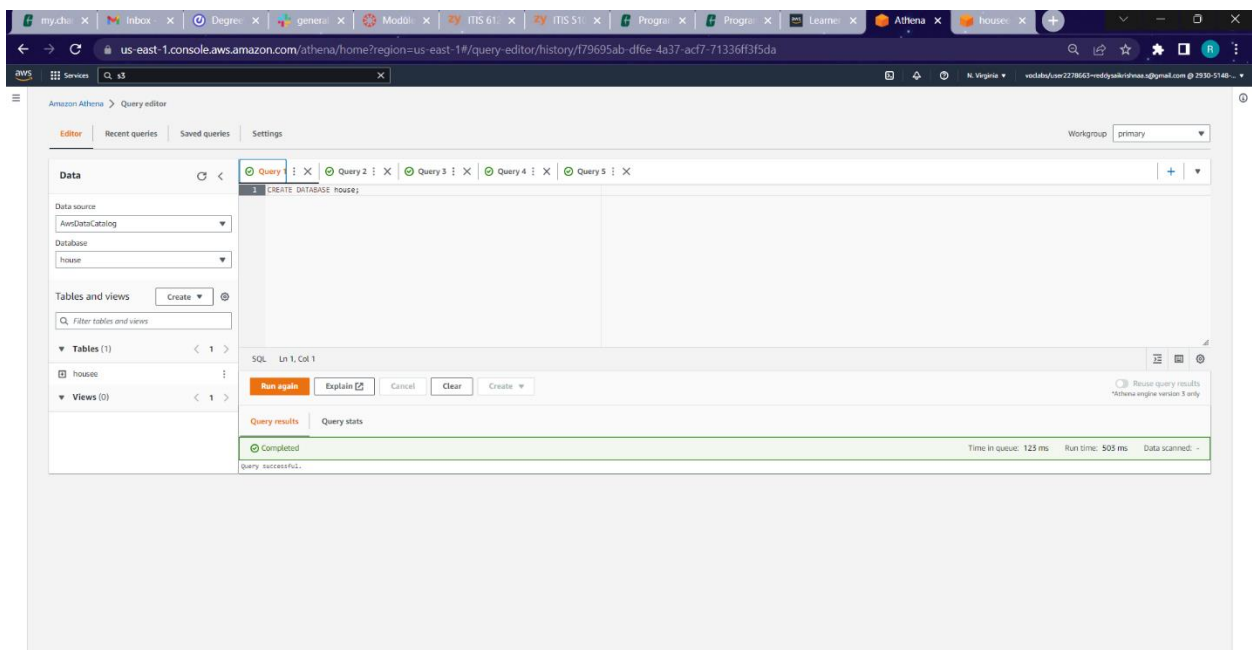
a) Exploratory Data Analysis

### AWS Components Used:

Amazon S3 Bucket

<https://housee.s3.amazonaws.com/housing.csv>

### Amazon Athena:



Amazon Athena > Query editor

Editor Recent queries Saved queries Settings

Workgroup: primary

Data source: AwsDataCatalog Database: house

Tables and views: Create Filter tables and views

Tables (1): house Views (0)

Query 1: X Query 2: X Query 3: X Query 4: X Query 5: X

1 SELECT \* FROM "house"."houseee" limit 10;

SQL Ln 1, Col 41

Run again Explain Cancel Clear Create

Query results Query stats

Completed Time in queue: 122 ms Run time: 576 ms Data scanned: 630.94 KB

Results (10)

Search rows

#	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_values	ocean_proximity
10	-122.25	37.84	42	2555	665	1206	595	2.0804	226700	NEAR BAY
9	-122.25	37.84	52	3104	687	1157	647	5.12	241400	NEAR BAY
8	-122.25	37.84	52	2535	489	1094	514	3.6591	299200	NEAR BAY
7	-122.25	37.85	52	919	213	413	193	4.0368	269700	NEAR BAY
6	-122.25	37.85	52	1627	280	565	259	3.8462	342200	NEAR BAY

Amazon Athena > Query editor

Editor Recent queries Saved queries Settings

Workgroup: primary

Data source: AwsDataCatalog Database: house

Tables and views: Create Filter tables and views

Tables (1): house Views (0)

Query 1: X Query 2: X Query 3: X Query 4: X Query 5: X

1 select housing\_median\_age, max(median\_house\_values) as median\_value from houseee group by housing\_median\_age;

SQL Ln 1, Col 109

Run again Explain Cancel Clear Create

Query results Query stats

Completed Time in queue: 160 ms Run time: 1.005 sec Data scanned: 1.14 MB

Results (53)

Search rows

#	housing_median_age	median_value
1	51	500001
2	46	500001
3	17	500001
4	25	500001
5	10	500001

The screenshot shows the Amazon Athena Query Editor interface. The top navigation bar includes tabs for 'my.ch...', 'Inbox', 'Degree', 'gener...', 'Module', 'ITIS 61', 'ITIS 51', 'Progr...', 'Progr...', 'Learn...', 'Athena', and 'house...'. The browser address bar shows the URL: `us-east-1.console.aws.amazon.com/athena/home?region=us-east-1#/query-editor/history/f699c548-d0b0-4ae1-81ff-37d4e48f4e41`.

The interface is divided into several sections:

- Editor:** Contains a SQL query: `1 SELECT sum(total_rooms) as rooms_sum FROM house;`
- Data:** Includes a 'Data source' dropdown set to 'AwsDataCatalog', a 'Database' dropdown set to 'house', and a 'Tables and views' section with a 'Create' button and a search bar.
- Tables (1):** Lists the 'house' table.
- Views (0):** No views are listed.
- Query results:** Shows the query status as 'Completed'. It includes a table with the following data:
 

#	rooms_sum
1	54402150

## Amazon Glue:

The screenshot shows the Amazon Glue Crawlers console. The top navigation bar includes tabs for 'my.e...', 'Inbox', 'Degree', 'Quer...', 'AWS', 'AWS', 'A X', 'AWS', 'Hon...', 'hou...', 'gen...', 'Cale...', 'Cale...', 'ITIS', 'ITIS', 'Progr...', 'Progr...', 'ame...', 'Step', and '+'. The browser address bar shows the URL: `us-east-1.console.aws.amazon.com/glue/home?region=us-east-1#/v2/data-catalog/crawlers/`.

The interface includes a left-hand navigation menu with the following sections:

- Data Catalog:** Databases, Tables, Stream schema registries, Schemas, Connections, Crawlers, Classifiers, Catalog settings.
- Data Integration and ETL:** AWS Glue Studio, Jobs, Interactive Sessions, Notebooks, Data classification tools, Sensitive data detection, Record Matching, Triggers, Workflows, Blueprints, Security configurations.
- Legacy pages:** What's New.

The main content area displays the 'Crawlers' section. It includes a notification banner about the new AWS Glue Crawlers console experience. Below the banner, there is a 'Crawlers (1) info' section with a search bar and a table of crawlers.

Name	State	Schedule	Last run	Log	Table changes from last run
houseee	Ready	-	-	-	-

us-east-1.console.aws.amazon.com/glue/home?region=us-east-1#/v2/data-catalog/tables/view/housing\_csv?database=house&catalogId=258066616206&versionId=latest

**Table details** | Advanced properties

Name housing_csv	Description -	Database house	Classification csv
Location s3://house/housing_csv	Connection -	Deprecated -	Last updated November 16, 2022 at 20:14:40
Input format org.apache.hadoop.mapred.TextInputFormat	Output format org.apache.hadoop.hive.gl.is.HiveIgnoreKeyTextOutputFormat	Serialize serialization lib org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe	

**Schema** | Partitions | Indexes

Schema (10)  
View and manage the table schema.

Filter schemas

#	Column name	Data type	Partition key	Comment
1	longitude	double	-	-
2	latitude	double	-	-
3	housing_median_age	bigint	-	-
4	total_rooms	bigint	-	-
5	total_bedrooms	bigint	-	-
6	population	bigint	-	-
7	households	bigint	-	-
8	median_income	double	-	-
9	median_house_value	bigint	-	-
10	ocean_proximity	string	-	-

## Amazon SageMaker:

### Dashboards:

house-dsglnotebook.us-east-1.sagemaker.aws/notebooks/housepriceprediction.ipynb

Jupyter | housepriceprediction | Last Checkpoint: a minute ago (unsaved changes) | Logout

File | Edit | View | Insert | Cell | Kernel | Widgets | Help | Not Trusted | conda\_python3

```
In [2]: #importing dataset into pandas as dataframe
DataFrame = pd.read_csv('s3://house/housing.csv')
ds=DataFrame #defining ds as a dataframe

In [3]: ds

Out[3]:
```

	longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income	median_house_value	ocean_proximity
0	-122.23	37.88	41	890	129.0	322	126	8.3252	452600	NEAR BAY
1	-122.22	37.86	21	7099	1106.0	2401	1138	8.3014	358500	NEAR BAY
2	-122.24	37.85	52	1467	190.0	496	177	7.2574	352100	NEAR BAY
3	-122.25	37.85	52	1274	235.0	558	219	5.6431	341300	NEAR BAY
4	-122.25	37.85	52	1627	280.0	565	259	3.8462	342200	NEAR BAY
...	...	...	...	...	...	...	...	...	...	...
20635	-121.09	39.48	25	1665	374.0	845	330	1.5603	78100	INLAND
20636	-121.21	39.49	18	697	150.0	356	114	2.5568	77100	INLAND
20637	-121.22	39.43	17	2254	485.0	1007	433	1.7000	92300	INLAND
20638	-121.32	39.43	18	1860	409.0	741	349	1.8672	84700	INLAND
20639	-121.24	39.37	16	2785	616.0	1387	530	2.3886	89400	INLAND

20640 rows x 10 columns

```
In [4]: #displays information of the dataset
ds.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20640 entries, 0 to 20639
Data columns (total 10 columns):
#   Column              Non-Null Count  Dtype
---  --
0   longitude            20640 non-null  float64
```

