

PROJECT REPORT

SPRING 2023 CPSC 531-03

HOUSE RENTAL ANALYSIS

TEAM MEMBERS

NAME :MANIKRISHNA SANGANABATLA
EMAIL: msanganabatla@csu.fullerton.edu
CWID:885638767

NAME: KRISHNA ABHISHITH PURIHELLA
EMAIL: purihella.abhi@csu.fullerton.edu
CWID:885211615

CONTENTS

1. INTRODUCTION
2. PROJECT OBJECTIVES
3. SCOPE AND LIMITATIONS
4. FUNCTIONALITIES
5. ARCHITECTURE AND DESIGN
6. GITHUB LOCATION OF CODE
7. DEPLOYMENT INSTRUCTIONS
8. STEPS TO RUN THE APPLICATION
9. TEST RESULTS.
10. REFERENCES.

INTRODUCTION:

This Project is centered on analyzing various aspects linked to a house. Prediction of house rent is vital among those aspects. Let us deep dive into some attributes that link us in analyzing the house rent. The price of a house depends on the region, area of the house, the facilities, etc., and prediction is carried out by using some of the existing machine learning algorithms. This Report focuses on the above part and how the data is collected and driven to do the analysis. Since this project requires a vast amount of data to analyze, the data is retrieved from the cloud and is then subjected to pre-processing and partitioning using AWS (Amazon Web Service) EMR clusters and is then analyzed by a tool (Athena). The analysis is then visualized by integrating Athena to one of the significant visualization tools called Tableau.

PROJECT OBJECTIVES:

GENERAL OBJECTIVE: The general objective of this project is to develop a visualization view of the rent that was predicted using various aspects.

SPECIFIC OBJECTIVE: Some of the specific objectives of this project are

- 1) Pre-process the dataset with several defects in the data—partitioning, and analyze the data by creating an EMR cluster that runs on AWS EC2 instances.
- 2) Integration of cluster instances with the system through SSH (Secure Socket Shell). A total of 3 instances are created – 1 Master Node and 2 Worker Nodes.
- 3) Create a PySpark program to query various scenarios associated with the dataset.
- 4) Submit a spark job and run the program
- 5) Upload the queried results into S3 bucket of AWS, where it can be again used to generate visualization charts.

- 6) Predict the House Rents using the Random-Forest method; here, we use a machine learning algorithm to predict various aspects using the dataset.
- 7) Application of big queries to generate tables of required attributes using AWS Athena.
- 8) Integration of AWS Athena with a visualization tool called Tableau.

SCOPE AND LIMITATIONS:

The scope of a house rental analysis project can differ depending on the project's objectives, accessible data, and anticipated findings. Examining rental costs can reveal the typical rentals charged for various types of homes in a given area. It can also detect trends over time and assist landlords, property managers, and renters make informed decisions.

Limitations of this project can be

1. Data Availability: One of the most significant restrictions of a house rental study project is data availability and quality. Rental data may be difficult to collect or incomplete, affecting the analysis's accuracy.
2. Lack of Standardization: Another constraint of house rental analysis is the lack of market standardization. Each market may have its distinct traits and elements influencing rental costs, making it difficult to compare statistics across regions.
3. The sample size of rental data may be too small to produce valid results or too vast to handle properly, limiting the scope of a rental analysis project.

FUNCTIONALITIES:

1)DATASET :

The Dataset is obtained from one of the open-source websites called Kaggle. Our dataset is of 1.2 Gb in total size. It is comprised of different attributes. There are a total of 22 columns like price, region, URL, sqfeet, type, region, state, etc.,

id	url	region	region_uri	price	type	sqfeet	beds	baths	cats_allow	dogs_allow	smoking	wheelchair	electric	vcomes_full	laundry	o_parking	o_image_url	description	lat	long	state
2	7E+09 https://bh.birmingha...	1195	apartment	1908	3	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
3	7E+09 https://bh.birmingha...	1120	apartment	1319	3	2	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
4	7E+09 https://bh.birmingha...	825	apartment	1183	1	1.5	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
5	7E+09 https://bh.birmingha...	800	apartment	927	1	1	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
6	7E+09 https://bh.birmingha...	785	apartment	1047	2	1	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
7	7E+09 https://bh.birmingha...	900	apartment	1298	2	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
8	7E+09 https://bh.birmingha...	925	apartment	1350	2	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
9	7E+09 https://bh.birmingha...	1355	apartment	1375	3	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4283	-86.7183	al
10	7E+09 https://bh.birmingha...	1120	apartment	1319	3	2	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
11	7E+09 https://bh.birmingha...	685	house	672	2	1	1	1	1	0	0	0	0	0	0	0	0	0	33.5755	-86.8856	al
12	7E+09 https://bh.birmingha...	1060	apartment	1100	2	2	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
13	7E+09 https://bh.birmingha...	1100	apartment	1100	2	2	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
14	7E+09 https://bh.birmingha...	870	apartment	851	1	1	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
15	7E+09 https://bh.birmingha...	815	apartment	851	1	1	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
16	7E+09 https://bh.birmingha...	1195	apartment	1908	3	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
17	7E+09 https://bh.birmingha...	1132	apartment	1107	2	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4816	-86.859	al
18	7E+09 https://bh.birmingha...	1105	apartment	1908	3	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
19	7E+09 https://bh.birmingha...	820	apartment	851	1	1	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
20	7E+09 https://bh.birmingha...	925	apartment	1350	2	2	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
21	7E+09 https://bh.birmingha...	785	apartment	1047	2	1	1	1	1	1	0	0	0	0	0	0	0	0	33.4226	-86.7065	al
22	7E+09 https://bh.birmingha...	890	apartment	1100	2	2	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
23	7E+09 https://bh.birmingha...	1120	apartment	1319	3	2	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al
24	7E+09 https://bh.birmingha...	975	apartment	1009	2	1	1	1	1	1	0	0	0	0	0	0	0	0	33.3755	-86.8045	al

2) Data Cleaning:

Data cleaning is a critical process to ensure data is correct and consistent. It involves removing duplicates, fixing errors, and removing irrelevant data. To do this, PySpark was installed, created a spark session, and created a custom dataset as a data frame. To clean the data, PySpark used commands like dropna() andfillna() to replace null values, handle invalid rows and columns, and replace null values with 0's .dropna() used to remove invalid/ unnecessary rows and columns and fillna() used to fill those rows/columns with the custom values.

AWS EMR Cluster: Amazon EMR (Elastic MapReduce) is a web service provided by Amazon Web Services (AWS) that allows customers to process massive volumes of data in a distributed computing environment. An EMR cluster is a collection of Amazon Elastic Compute Cloud (EC2) machines running Hadoop and other large data processing frameworks collaborating to accomplish data processing tasks. Here are some of the important features of an AWS EMR cluster:

Distributed Processing: The fundamental role of an EMR cluster is to enable distributed processing of huge datasets. This means data can be broken into smaller bits and processed in parallel across numerous EC2 instances.

EMR clusters support a variety of large data processing frameworks, including Apache Hadoop, Spark, and Hive.

Scalability: EMR clusters are extremely scalable, allowing customers to add or delete EC2 instances as needed to fulfill processing requirements. This makes it simple to scale up or down based on changing data processing requirements.

Automation: EMR clusters automate many of the procedures associated with constructing and operating Hadoop clusters, such as installation, configuration, and scaling.

Integration with other AWS Services: EMR clusters interface with other AWS services such as S3, EC2, and CloudWatch, allowing users to move data effortlessly across services and monitor cluster performance.

APACHE SPARK: Apache Spark is a distributed computing system that is open source and used for big data processing. It has many features to help with data processing, analysis, and machine learning activities. The key functionalities are In-Memory Computing, Distributed Computing, Data Processing, Machine Learning, and Integration with Other Tools like Hadoop, Cassandra, and Amazon S3, making working with data from multiple sources easy.

AWS S3 BUCKET: Amazon Simple Storage Service (S3) is an object storage service that is extremely scalable and provided by Amazon Web Services (AWS). It offers a wide range of features to aid in the storage, management, and access of data objects in the cloud. Here are some of the key features of an AWS S3 bucket:

Scalability: S3 offers limitless scalability, allowing customers to store and retrieve massive volumes of data objects in the cloud.

Storage: S3 provides a durable storage solution designed to withstand data loss due to hardware failures, natural disasters, and other occurrences.

Availability: S3 delivers high availability for data objects by storing several redundant copies of data across different availability zones.

Security: S3 offers a variety of security capabilities to safeguard data objects, including encryption of data at rest and in transit, access control via AWS Identity and Access Management (IAM) policies, and versioning to keep track of changes to data objects.

Integration with Other AWS Services: S3 integrates with other AWS services, such as EC2, Lambda, and EMR, enabling users to easily store and access data objects from different services.

AWS Athena:

Amazon Athena is an interactive query service that allows users to run conventional SQL queries to explore data in Amazon S3. It has a wide range of functions to help with data analysis and reporting. Some of the important aspects are as follows:

Scalability: Because Athena is highly scalable and capable of handling enormous datasets, users can examine data in S3 at any scale.

Compatibility: Because Athena supports normal SQL queries, customers can easily analyze their data using familiar SQL syntax. It also supports various data formats such as CSV, JSON, ORC, Parquet, and Avro.

Athena interfaces with other AWS services, such as AWS Glue and Amazon QuickSight, allowing customers to easily extract, manipulate, and visualize their data.

Performance: Depending on the complexity of the query, Athena provides quick query performance, with results given in seconds or minutes.

TABLEAU: Tableau is a powerful data visualization and business intelligence software that businesses and organizations use to analyze, interpret, and exchange data. Tableau's primary features are as follows:

Data Connectivity: Tableau can link to a variety of data sources, including spreadsheets, databases, and cloud-based data warehouses. It also supports several file types, including CSV, PARQUET, Excel, and JSON.

Data Preparation: Tableau provides strong data preparation features that enable users to transform, clean, and blend data from many sources to generate a cohesive dataset. This allows users to arrange their data for analysis and visualization without the need for complicated code.

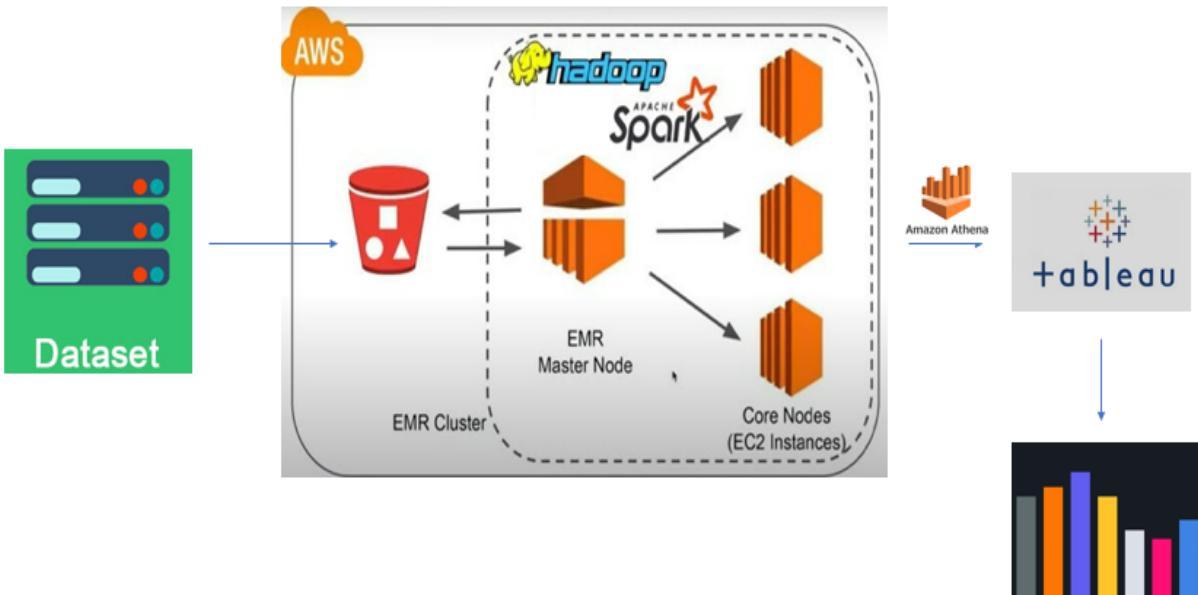
Data Visualization: Tableau provides a wide range of visualization tools to create interactive and engaging visualizations like charts, graphs, and maps. Users can also create dashboards and stories to share their insights with others.

These are some of the detailed functionalities of the tools and services used to design this project.

APPROACH: Below is the following approach for our project :



DESIGN AND ARCHITECTURE:



GITHUB LOCATION OF CODE :

https://github.com/manikittu810/ADBMS_PROJECT

TECHNOLOGY USED :

OS: Windows , macOS

Framework : Amazon Web Service (AWS) , Spark , for BigQuery used SQL

Languages Used : Python, SQL

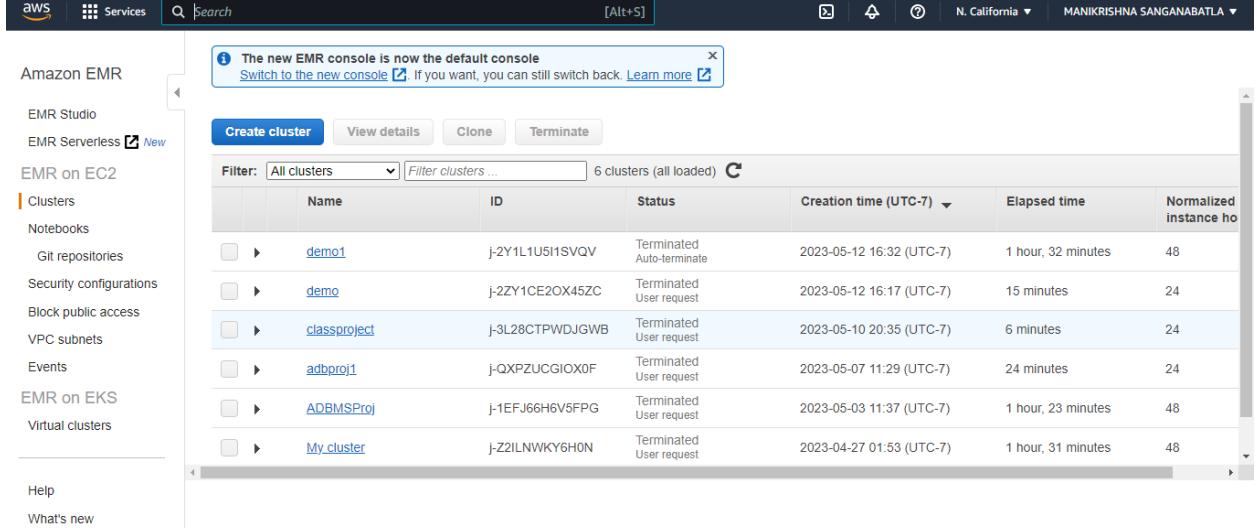
Libraries : Pyspark

IDE : Jupyter Notebook ,VS code

Below is a design diagram of this project that represents the flow and usage of the tools and applications.

DESIGN :

- A cleaned Dataset is taken and is uploaded into the S3 Bucket of the AWS. This S3 bucket is the cloud data storage provided by AWS.

- Then comes the creation of EMR Clusters(Elastic Map Reduce Cluster). EMR cluster is a hadoop cluster that helps in generating instances to run the data on distributed file systems.
- The EMR cluster is generated by giving necessary parameters, here we have used the Apache Hadoop yarn version and this cluster consists of total 3 instances (1 master and 2 slave nodes). **A cluster with the name of <any_custom_name> is created in this project.**
- The screenshot shows the AWS EMR console interface. On the left, there's a sidebar with navigation links: Amazon EMR, EMR Studio, EMR Serverless (New), EMR on EC2 (selected), Clusters, Notebooks, Git repositories, Security configurations, Block public access, VPC subnets, Events, EMR on EKS, and Virtual clusters. At the top, there's a search bar and a message: "The new EMR console is now the default console. Switch to the new console If you want, you can still switch back. Learn more". Below the message are buttons for "Create cluster", "View details", "Clone", and "Terminate". A table lists six clusters: demo1, demo, classproject, adbproj1, ADBMSProj, and My_cluster. Each row includes columns for Name, ID, Status, Creation time (UTC-7), Elapsed time, and Normalized instance hours.

- A python code is designed locally that queries various aspects; the data stored in the S3 bucket is retrieved in this code to provide the output. You can find this code inside the GitHub link provided in this document with the **Filename <main.py> and <main1.py>**

```
File Edit Selection View Go Run Terminal Help main1.py - Visual Studio Code
Restricted Mode is intended for safe code browsing. Trust this window to enable all features. Manage Learn More

main1.py 2 × main.py 2
C: > Users > msanganabatla > Desktop > proj > main1.py > main
1   from pyspark.sql import SparkSession
2   from pyspark.sql.functions import col
3
4   S3_DATA_SOURCE_PATH = 's3://classnow/class/houseRent.csv'
5   S3_DATA_OUTPUT_PATH = 's3://classnow/in_class_output'
6
7
8 def main():
9
10    spark = SparkSession.builder.appName('ADBMSPProject').getOrCreate()
11
12    all_data = spark.read.csv(S3_DATA_SOURCE_PATH,header=True)
13
14    print('Total number of records in the source data :%s' % all_data.count())
15
16
17    selected_data = all_data.where((col('beds') == 2) & (col('sqfeet')>1500))
18
19    print('Total number of beds with 2 are :%s'%selected_data.count())
20
21
22    selected_data = all_data.where((col('beds') == 1) & (col('sqfeet')>1500))
23
24    print('Total number of beds with 1 are :%s'%selected_data.count())
25
26    selected_data = all_data.where((col('baths')==1) & (col('type')=='apartment'))
27
28
29
30
31
```

You can find the same code in the github using the above provided link.

- Integration of the EMR cluster with the local using SSH is carried out. This is done submitting the following commands on the command prompt.
 - `Ssh -i <location (or) path of the aws key pair> <cluster provided command for integrating master node with local>`

- A spark job is then submitted on the local console by using the command “**spark-submit <filename.py>**”. After submission of the spark job, the code is then run and the results are displayed on the console.

```
[hadoop@ip-172-31-16-97 ~]$ vi main.py
[hadoop@ip-172-31-16-97 ~]$ spark-submit main.py
23/05/12 23:39:55 INFO SparkContext: Running Spark version 2.4.8-amzn-2
23/05/12 23:39:55 INFO SparkContext: Submitted application: ADBMSProject
23/05/12 23:39:55 INFO SecurityManager: Changing view acls to: hadoop
23/05/12 23:39:55 INFO SecurityManager: Changing modify acls to: hadoop
23/05/12 23:39:55 INFO SecurityManager: Changing view acls groups to:
23/05/12 23:39:55 INFO SecurityManager: Changing modify acls groups to:
23/05/12 23:39:55 INFO SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(hadoop); groups with view permissions: Set(); users with modify permissions: Set(hadoop); groups with modify permissions: Set()
23/05/12 23:39:55 INFO Utils: Successfully started service 'sparkDriver' on port 45761.
23/05/12 23:39:55 INFO SparkEnv: Registering MapOutputTracker
23/05/12 23:39:55 INFO SparkEnv: Registering BlockManagerMaster
23/05/12 23:39:55 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
23/05/12 23:39:55 INFO BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
23/05/12 23:39:55 INFO DiskBlockManager: Created local directory at /mnt/tmp/blockmgr-a9154bd6-8f8d-4172-846b-eb07c814c12b
23/05/12 23:39:55 INFO MemoryStore: MemoryStore started with capacity 912.3 MB
23/05/12 23:39:55 INFO SparkEnv: Registering OutputCommitCoordinator
23/05/12 23:39:56 INFO SparkUI: Successfully started service 'SparkUI' on port 4048.
23/05/12 23:39:56 INFO SparkUI: Bound SparkUI to 0.0.0.0, and started at http://ip-172-31-16-97.us-west-1.compute.internal:4048
23/05/12 23:39:56 INFO Utils: Using 50 preallocated executors (minExecutors: 0). Set spark.dynamicAllocation.preallocateExecutors to 'false' disable executor preallocation.
23/05/12 23:39:57 INFO RMPProxy: Connecting to ResourceManager at ip-172-31-16-97.us-west-1.compute.internal/172.31.16.97:8032
23/05/12 23:39:57 INFO Client: Requesting a new application from cluster with 2 NodeManagers
23/05/12 23:39:57 INFO Configuration: resource-types.xml not found
23/05/12 23:39:57 INFO ResourceUtils: Unable to find 'resource-types.xml'.
23/05/12 23:39:57 INFO ResourceUtils: Adding resource type - name = memory_mb, units = Mi, type = COUNTABLE
23/05/12 23:39:57 INFO ResourceUtils: Adding resource type - name = vcores, units = , type = COUNTABLE
23/05/12 23:39:57 INFO Client: Verifying our application has not requested more than the maximum memory capability of the cluster (12288 MB per container)
23/05/12 23:39:57 INFO Client: Will allocate AM container, with 896 MB memory including 384 MB overhead
23/05/12 23:39:57 INFO Client: Setting up container launch context for our AM
23/05/12 23:39:57 INFO Client: Setting up the launch environment for our AM container
23/05/12 23:39:57 INFO Client: Preparing resources for our AM container
23/05/12 23:39:57 WARN Client: Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
```

- The output results are stored back to S3 bucket into a different folder , whose path is already defined in the code that was ran by submitting the spark job. These files are again used in the later stages to generate visual reports.

Buckets (5) Info
Buckets are containers for data stored in S3. [Learn more](#)

Name	AWS Region	Access	Creation date
athenaproj	US West (N. California) us-west-1	Bucket and objects not public	May 12, 2023, 17:39:12 (UTC-07:00)
aws-logs-337158828970-us-east-1	US East (N. Virginia) us-east-1	Bucket and objects not public	May 12, 2023, 16:10:37 (UTC-07:00)
aws-logs-337158828970-us-west-1	US West (N. California) us-west-1	Bucket and objects not public	May 12, 2023, 16:17:29 (UTC-07:00)
classnow	US East (N. Virginia) us-east-1	Bucket and objects not public	May 10, 2023, 23:57:50 (UTC-07:00)
		Bucket	

classnow Info

- [Objects](#) Objects
- [Properties](#)
- [Permissions](#)
- [Metrics](#)
- [Management](#)
- [Access Points](#)

Objects (4)

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](#)

Name	Type	Last modified	Size	Storage class
class/	Folder	-	-	-
in_class_output_\$folder\$	-	May 12, 2023, 16:43:43 (UTC-07:00)	0 B	Standard
in_class_output/	Folder	-	-	-
in_class_output1/	Folder	-	-	-

ML Algorithm for prediction : For predictions we have used random forest algorithm which gave us the accurate results. We have also used many algorithms but random forest gave us the best results. We have many attributes in our dataset , so, we have only taken four attributes to predict the rent of a type of a house in the state of california. Below is the code and results.

```

# Import standard libraries
import os
import pandas as pd
from sklearn.model_selection import train_test_split, RepeatedKFold
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error
from xgboost import XGBRegressor

# Read the data
df = pd.read_csv('houseprice.csv')

# Split the data into features (X) and target variable (y)
X = df.drop(['SalePrice'], axis=1).values
y = df['SalePrice'].values

# One-hot encode categorical features
X_cat = pd.get_dummies(X[:, 3], columns=['MSSubClass'])
X_cat = pd.get_dummies(X_cat[:, 5], columns=['OverallQual', 'OverallCond'])
X_cat = pd.get_dummies(X_cat[:, 6], columns=['ExterQual', 'ExterCond'])

scaler = StandardScaler()
X = pd.concat([X_cat, X[:, 7:-1], pd.DataFrame(scaler.fit_transform(X[:, -1]), columns=X[:, -1])], axis=1)

# Create the categorical and numerical features
X = pd.concat([X_cat, X[:, -1], X[:, -1]], axis=1)

# Create the training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Fine-tune the model using RepeatedKFold
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [5, 10, 15],
    'learning_rate': [0.01, 0.05, 0.1],
    'min_child_weight': [1, 5, 10],
    'subsample': [0.8, 0.9, 1.0]
}

# Create the RepeatedKFold object
repeated_search = RepeatedKFold(n_splits=5, n_repeats=3, random_state=42)

# Fit the best parameters and score from the repeat search
for (train_index, test_index) in repeated_search.split(X_train, y_train):
    print('Next Score:', repeated_search.score_)

# Fit the model with the best hyperparameters
xgb = XGBRegressor(**repeated_search.get_params())
xgb.fit(X_train, y_train)

# Make predictions on the test set
xgb_pred = xgb.predict(X_test)

# Print the RMSE
print(f'RMSE: {mean_squared_error(y_test, xgb_pred)}')

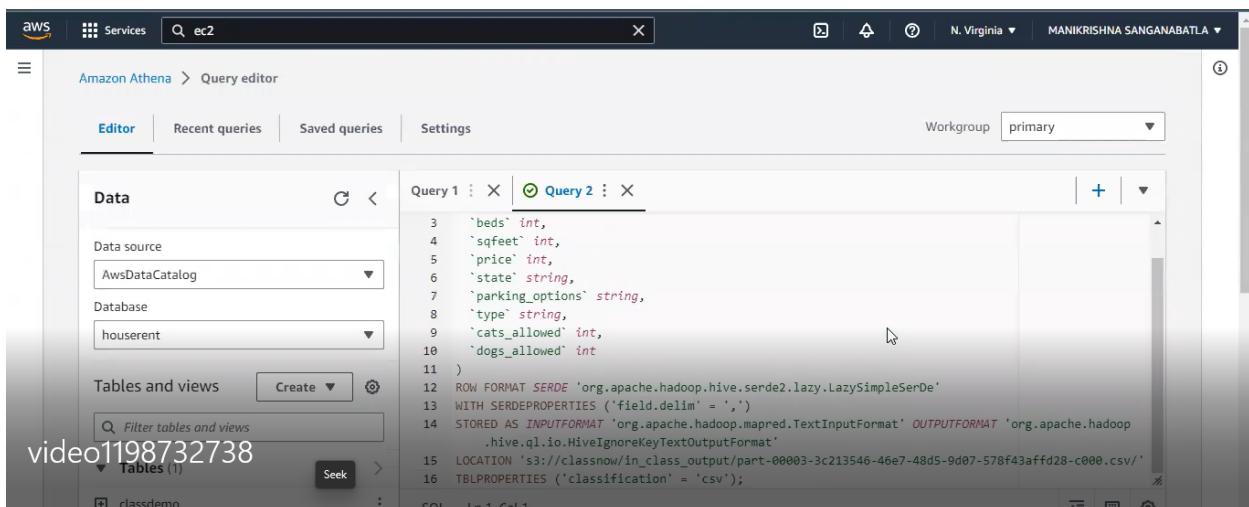
# Save the data
new_data = pd.DataFrame([
    {'id': 1, 'MSSubClass': 2019, 'OverallQual': 3, 'OverallCond': 3, 'YearBuilt': 1890, 'YearRemodAdd': 1890, 'TotalBsmtSF': 0, '1stFlrSF': 500, '2ndFlrSF': 0, 'LowQualFinSF': 0, 'GrLivArea': 850, 'BsmtUnfSF': 0, 'BsmtFinSF1': 0, 'BsmtFinSF2': 0, 'BsmtFullBath': 0, 'BsmtHalfBath': 0, 'FullBath': 1, 'HalfBath': 0, 'BedroomAbvGr': 1, 'KitchenAbvGr': 1, 'TotRmsAbvGrd': 3, 'FireplaceQu': 'None', 'GarageCars': 0, 'GarageArea': 0, 'GarageYrBlt': 0, 'GarageCond': 'None', 'PavedDrive': 'N', 'WoodDeckSF': 0, 'OpenPorchSF': 0, 'EnclosedPorch': 0, 'ScreenPorch': 0, 'PoolArea': 0, 'PoolQC': 'None', 'Fence': 'None', 'MiscVal': 0, 'MoSold': 0, 'YrSold': 1978, 'SaleType': 'Normal', 'SaleCondition': 'Normal'}])
new_data.to_csv('new_data.csv', index=False)

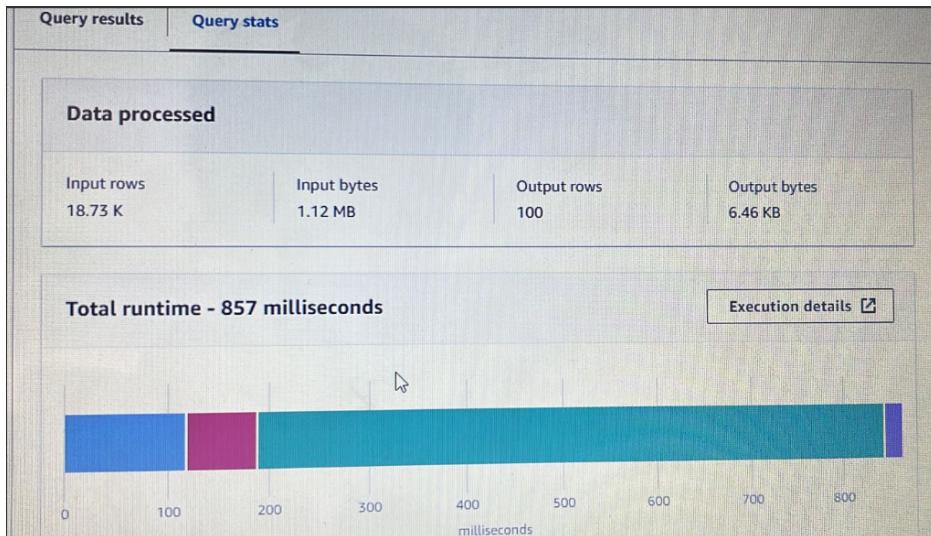
my_data = pd.read_csv('new_data.csv', index_col='id')
my_data['MSSubClass'] = my_data['MSSubClass'].map({2019: 1, 3120: 2, 3121: 3, 3122: 4, 3123: 5, 3124: 6, 3125: 7, 3126: 8, 3127: 9, 3128: 10, 3129: 11, 3130: 12, 3131: 13, 3132: 14, 3133: 15, 3134: 16, 3135: 17, 3136: 18, 3137: 19, 3138: 20, 3139: 21, 3140: 22, 3141: 23, 3142: 24, 3143: 25, 3144: 26, 3145: 27, 3146: 28, 3147: 29, 3148: 30, 3149: 31, 3150: 32, 3151: 33, 3152: 34, 3153: 35, 3154: 36, 3155: 37, 3156: 38, 3157: 39, 3158: 40, 3159: 41, 3160: 42, 3161: 43, 3162: 44, 3163: 45, 3164: 46, 3165: 47, 3166: 48, 3167: 49, 3168: 50, 3169: 51, 3170: 52, 3171: 53, 3172: 54, 3173: 55, 3174: 56, 3175: 57, 3176: 58, 3177: 59, 3178: 60, 3179: 61, 3180: 62, 3181: 63, 3182: 64, 3183: 65, 3184: 66, 3185: 67, 3186: 68, 3187: 69, 3188: 70, 3189: 71, 3190: 72, 3191: 73, 3192: 74, 3193: 75, 3194: 76, 3195: 77, 3196: 78, 3197: 79, 3198: 80, 3199: 81, 3200: 82, 3201: 83, 3202: 84, 3203: 85, 3204: 86, 3205: 87, 3206: 88, 3207: 89, 3208: 90, 3209: 91, 3210: 92, 3211: 93, 3212: 94, 3213: 95, 3214: 96, 3215: 97, 3216: 98, 3217: 99, 3218: 100, 3219: 101, 3220: 102, 3221: 103, 3222: 104, 3223: 105, 3224: 106, 3225: 107, 3226: 108, 3227: 109, 3228: 110, 3229: 111, 3230: 112, 3231: 113, 3232: 114, 3233: 115, 3234: 116, 3235: 117, 3236: 118, 3237: 119, 3238: 120, 3239: 121, 3240: 122, 3241: 123, 3242: 124, 3243: 125, 3244: 126, 3245: 127, 3246: 128, 3247: 129, 3248: 130, 3249: 131, 3250: 132, 3251: 133, 3252: 134, 3253: 135, 3254: 136, 3255: 137, 3256: 138, 3257: 139, 3258: 140, 3259: 141, 3260: 142, 3261: 143, 3262: 144, 3263: 145, 3264: 146, 3265: 147, 3266: 148, 3267: 149, 3268: 150, 3269: 151, 3270: 152, 3271: 153, 3272: 154, 3273: 155, 3274: 156, 3275: 157, 3276: 158, 3277: 159, 3278: 160, 3279: 161, 3280: 162, 3281: 163, 3282: 164, 3283: 165, 3284: 166, 3285: 167, 3286: 168, 3287: 169, 3288: 170, 3289: 171, 3290: 172, 3291: 173, 3292: 174, 3293: 175, 3294: 176, 3295: 177, 3296: 178, 3297: 179, 3298: 180, 3299: 181, 3300: 182, 3301: 183, 3302: 184, 3303: 185, 3304: 186, 3305: 187, 3306: 188, 3307: 189, 3308: 190, 3309: 191, 3310: 192, 3311: 193, 3312: 194, 3313: 195, 3314: 196, 3315: 197, 3316: 198, 3317: 199, 3318: 200, 3319: 201, 3320: 202, 3321: 203, 3322: 204, 3323: 205, 3324: 206, 3325: 207, 3326: 208, 3327: 209, 3328: 210, 3329: 211, 3330: 212, 3331: 213, 3332: 214, 3333: 215, 3334: 216, 3335: 217, 3336: 218, 3337: 219, 3338: 220, 3339: 221, 3340: 222, 3341: 223, 3342: 224, 3343: 225, 3344: 226, 3345: 227, 3346: 228, 3347: 229, 3348: 230, 3349: 231, 3350: 232, 3351: 233, 3352: 234, 3353: 235, 3354: 236, 3355: 237, 3356: 238, 3357: 239, 3358: 240, 3359: 241, 3360: 242, 3361: 243, 3362: 244, 3363: 245, 3364: 246, 3365: 247, 3366: 248, 3367: 249, 3368: 250, 3369: 251, 3370: 252, 3371: 253, 3372: 254, 3373: 255, 3374: 256, 3375: 257, 3376: 258, 3377: 259, 3378: 260, 3379: 261, 3380: 262, 3381: 263, 3382: 264, 3383: 265, 3384: 266, 3385: 267, 3386: 268, 3387: 269, 3388: 270, 3389: 271, 3390: 272, 3391: 273, 3392: 274, 3393: 275, 3394: 276, 3395: 277, 3396: 278, 3397: 279, 3398: 280, 3399: 281, 3400: 282, 3401: 283, 3402: 284, 3403: 285, 3404: 286, 3405: 287, 3406: 288, 3407: 289, 3408: 290, 3409: 291, 3410: 292, 3411: 293, 3412: 294, 3413: 295, 3414: 296, 3415: 297, 3416: 298, 3417: 299, 3418: 300, 3419: 301, 3420: 302, 3421: 303, 3422: 304, 3423: 305, 3424: 306, 3425: 307, 3426: 308, 3427: 309, 3428: 310, 3429: 311, 3430: 312, 3431: 313, 3432: 314, 3433: 315, 3434: 316, 3435: 317, 3436: 318, 3437: 319, 3438: 320, 3439: 321, 3440: 322, 3441: 323, 3442: 324, 3443: 325, 3444: 326, 3445: 327, 3446: 328, 3447: 329, 3448: 330, 3449: 331, 3450: 332, 3451: 333, 3452: 334, 3453: 335, 3454: 336, 3455: 337, 3456: 338, 3457: 339, 3458: 340, 3459: 341, 3460: 342, 3461: 343, 3462: 344, 3463: 345, 3464: 346, 3465: 347, 3466: 348, 3467: 349, 3468: 350, 3469: 351, 3470: 352, 3471: 353, 3472: 354, 3473: 355, 3474: 356, 3475: 357, 3476: 358, 3477: 359, 3478: 360, 3479: 361, 3480: 362, 3481: 363, 3482: 364, 3483: 365, 3484: 366, 3485: 367, 3486: 368, 3487: 369, 3488: 370, 3489: 371, 3490: 372, 3491: 373, 3492: 374, 3493: 375, 3494: 376, 3495: 377, 3496: 378, 3497: 379, 3498: 380, 3499: 381, 3500: 382, 3501: 383, 3502: 384, 3503: 385, 3504: 386, 3505: 387, 3506: 388, 3507: 389, 3508: 390, 3509: 391, 3510: 392, 3511: 393, 3512: 394, 3513: 395, 3514: 396, 3515: 397, 3516: 398, 3517: 399, 3518: 400, 3519: 401, 3520: 402, 3521: 403, 3522: 404, 3523: 405, 3524: 406, 3525: 407, 3526: 408, 3527: 409, 3528: 410, 3529: 411, 3530: 412, 3531: 413, 3532: 414, 3533: 415, 3534: 416, 3535: 417, 3536: 418, 3537: 419, 3538: 420, 3539: 421, 3540: 422, 3541: 423, 3542: 424, 3543: 425, 3544: 426, 3545: 427, 3546: 428, 3547: 429, 3548: 430, 3549: 431, 3550: 432, 3551: 433, 3552: 434, 3553: 435, 3554: 436, 3555: 437, 3556: 438, 3557: 439, 3558: 440, 3559: 441, 3560: 442, 3561: 443, 3562: 444, 3563: 445, 3564: 446, 3565: 447, 3566: 448, 3567: 449, 3568: 450, 3569: 451, 3570: 452, 3571: 453, 3572: 454, 3573: 455, 3574: 456, 3575: 457, 3576: 458, 3577: 459, 3578: 460, 3579: 461, 3580: 462, 3581: 463, 3582: 464, 3583: 465, 3584: 466, 3585: 467, 3586: 468, 3587: 469, 3588: 470, 3589: 471, 3590: 472, 3591: 473, 3592: 474, 3593: 475, 3594: 476, 3595: 477, 3596: 478, 3597: 479, 3598: 480, 3599: 481, 3600: 482, 3601: 483, 3602: 484, 3603: 485, 3604: 486, 3605: 487, 3606: 488, 3607: 489, 3608: 490, 3609: 491, 3610: 492, 3611: 493, 3612: 494, 3613: 495, 3614: 496, 3615: 497, 3616: 498, 3617: 499, 3618: 500, 3619: 501, 3620: 502, 3621: 503, 3622: 504, 3623: 505, 3624: 506, 3625: 507, 3626: 508, 3627: 509, 3628: 510, 3629: 511, 3630: 512, 3631: 513, 3632: 514, 3633: 515, 3634: 516, 3635: 517, 3636: 518, 3637: 519, 3638: 520, 3639: 521, 3640: 522, 3641: 523, 3642: 524, 3643: 525, 3644: 526, 3645: 527, 3646: 528, 3647: 529, 3648: 530, 3649: 531, 3650: 532, 3651: 533, 3652: 534, 3653: 535, 3654: 536, 3655: 537, 3656: 538, 3657: 539, 3658: 540, 3659: 541, 3660: 542, 3661: 543, 3662: 544, 3663: 545, 3664: 546, 3665: 547, 3666: 548, 3667: 549, 3668: 550, 3669: 551, 3670: 552, 3671: 553, 3672: 554, 3673: 555, 3674: 556, 3675: 557, 3676: 558, 3677: 559, 3678: 560, 3679: 561, 3680: 562, 3681: 563, 3682: 564, 3683: 565, 3684: 566, 3685: 567, 3686: 568, 3687: 569, 3688: 570, 3689: 571, 3690: 572, 3691: 573, 3692: 574, 3693: 575, 3694: 576, 3695: 577, 3696: 578, 3697: 579, 3698: 580, 3699: 581, 3700: 582, 3701: 583, 3702: 584, 3703: 585, 3704: 586, 3705: 587, 3706: 588, 3707: 589, 3708: 590, 3709: 591, 3710: 592, 3711: 593, 3712: 594, 3713: 595, 3714: 596, 3715: 597, 3716: 598, 3717: 599, 3718: 600, 3719: 601, 3720: 602, 3721: 603, 3722: 604, 3723: 605, 3724: 606, 3725: 607, 3726: 608, 3727: 609, 3728: 610, 3729: 611, 3730: 612, 3731: 613, 3732: 614, 3733: 615, 3734: 616, 3735: 617, 3736: 618, 3737: 619, 3738: 620, 3739: 621, 3740: 622, 3741: 623, 3742: 624, 3743: 625, 3744: 626, 3745: 627, 3746: 628, 3747: 629, 3748: 630, 3749: 631, 3750: 632, 3751: 633, 3752: 634, 3753: 635, 3754: 636, 3755: 637, 3756: 638, 3757: 639, 3758: 640, 3759: 641, 3760: 642, 3761: 643, 3762: 644, 3763: 645, 3764: 646, 3765: 647, 3766: 648, 3767: 649, 3768: 650, 3769: 651, 3770: 652, 3771: 653, 3772: 654, 3773: 655, 3774: 656, 3775: 657, 3776: 658, 3777: 659, 3778: 660, 3779: 661, 3780: 662, 3781: 663, 3782: 664, 3783: 665, 3784: 666, 3785: 667, 3786: 668, 3787: 669, 3788: 670, 3789: 671, 3790: 672, 3791: 673, 3792: 674, 3793: 675, 3794: 676, 3795: 677, 3796: 678, 3797: 679, 3798: 680, 3799: 681, 3800: 682, 3801: 683, 3802: 684, 3803: 685, 3804: 686, 3805: 687, 3806: 688, 3807: 689, 3808: 690, 3809: 691, 3810: 692, 3811: 693, 3812: 694, 3813: 695, 3814: 696, 3815: 697, 3816: 698, 3817: 699, 3818: 700, 3819: 701, 3820: 702, 3821: 703, 3822: 704, 3823: 705, 3824: 706, 3825: 707, 3826: 708, 3827: 709, 3828: 710, 3829: 711, 3830: 712, 3831: 713, 3832: 714, 3833: 715, 3834: 716, 3835: 717, 3836: 718, 3837: 719, 3838: 720, 3839: 721, 3840: 722, 3841: 723, 3842: 724, 3843: 725, 3844: 726, 3845: 727, 3846: 728, 3847: 729, 3848: 730, 3849: 731, 3850: 732, 3851: 733, 3852: 734, 3853: 735, 3854: 736, 3855: 737, 3856: 738, 3857: 739, 3858: 740, 3859: 741, 3860: 742, 3861: 743, 3862: 744, 3863: 745, 3864: 746, 3865: 747, 3866: 748, 3867: 749, 3868: 750, 3869: 751, 3870: 752, 3871: 753, 3872: 754, 3873: 755, 3874: 756, 3875: 757, 3876: 758, 3877: 759, 3878: 760, 3879: 761, 3880: 762, 3881: 763, 3882: 764, 3883: 765, 3884: 766, 3885: 767, 3886: 768, 3887: 769, 3888: 770, 3889: 771, 3890: 772, 3891: 773, 3892: 774, 3893: 775, 3894: 776, 3895: 777, 3896: 778, 3897: 779, 3898: 780, 3899: 781, 3900: 782, 3901: 783, 3902: 784, 3903: 785, 3904: 786, 3905: 787, 3906: 788, 3907: 789, 3908: 790, 3909: 791, 3910: 792, 3911: 793, 3912: 794, 3913: 795, 3914: 796, 3915: 797, 3916: 798, 3917: 799, 3918: 800, 3919: 801, 3920: 802, 3921: 803, 3922: 804, 3923: 805, 3924: 806, 3925: 807, 3926: 808, 3927: 809, 3928: 810, 3929: 811, 3930: 812, 3931: 813, 3932: 814, 3933: 815, 3934: 816, 3935: 817, 3936: 818, 3937: 819, 3938: 820, 3939: 821, 3940: 822, 3941: 823, 3942: 824, 3943: 825, 3944: 826, 3945: 827, 3946: 828, 3947: 829, 3948: 830, 3949: 831, 3950: 832, 3951: 833, 3952: 834, 3953: 835, 3954: 836, 3955: 837, 3956: 838, 3957: 839, 3958: 840, 3959: 841, 3960: 842, 3961: 843, 3962: 844, 3963: 845, 3964: 846, 3965: 847, 3966: 848, 3967: 849, 3968: 850, 3969: 851, 3970: 852, 3971: 853, 3972: 854, 3973: 855, 3974: 856, 3975: 857, 3976: 858, 3977: 859, 3978: 860, 3979: 861, 3980: 862, 3981: 863, 3982: 864, 3983: 865, 3984: 866, 3985: 867, 3986: 868, 3987: 869, 3988: 870, 3989: 871, 3990: 872, 3991: 873, 3992: 874, 3993: 875, 3994: 876, 3995: 877, 3996: 878, 3997: 879, 3998: 880, 3999: 881, 4000: 882, 4001: 883, 4002: 884, 4003: 885, 4004: 886, 4005: 887, 4006: 888, 4007: 889, 4008: 890, 4009: 891, 4010: 892, 4011: 893, 4012: 894, 4013: 895, 4014: 896, 4015: 897, 4016: 898, 4017: 899, 4018: 900, 4019: 901, 4020: 902, 4021: 903, 4022: 904, 4023: 905, 4024: 906, 4025: 907, 4026: 908, 4027: 909, 4028: 910, 4029: 911, 4030: 912, 4031: 913, 4032: 914, 4033: 915, 4034: 916, 4035: 917, 4036: 918, 4037: 919, 4038: 920, 4039: 921, 4040: 922, 4041: 923, 4042: 924, 4043: 925, 4044: 926, 4045: 927, 4046: 928, 4047: 929, 4048: 930, 4049: 931, 4050: 932, 4051: 933, 4052: 934, 4053: 935, 4054: 936, 4055: 937, 4056: 938, 4057: 939, 4058: 940, 4059: 941, 4060: 942, 4061: 943, 4062: 944, 4063: 945, 4064: 946, 4065: 947, 4066: 948, 4067: 949, 4068: 950, 4069: 951, 4070: 952, 4071: 953, 4072: 954, 4073: 955, 4074: 956, 4075: 957, 4076: 958, 4077: 959, 4078: 960, 4079: 961, 4080: 962, 4081: 963, 4082: 964, 4083: 965, 4084: 966, 4085: 967, 4086: 968, 4087: 969, 4088: 970, 4089: 971, 4090: 972, 4091: 973, 4092: 974, 4093: 975, 4094: 976, 4095: 977, 4096: 978, 4097: 979, 4098: 980, 4099: 981, 4100: 982, 4101: 983, 4102: 984, 4103: 985, 4104: 986, 4105: 987, 4106: 988, 4107: 989, 4108: 990, 4109: 991, 4110: 992, 4111: 993, 4112: 994, 4113: 995, 4114: 996, 4115: 997, 4116: 998, 4117: 999, 4118: 1000, 4119: 1001, 4120: 1002, 4121: 1003, 4122: 1004, 4123: 1005, 4124: 1006, 4125: 1007, 4126: 1008, 4127: 1009, 4128: 1010, 4129: 1011, 4130: 1012, 4131: 1013, 4132: 1014, 4133: 1015, 4134: 1016, 4135: 1017, 4136: 1018, 4137: 1019, 4138: 1020, 4139: 1021, 4140: 1022, 4141: 1023, 4142: 1024, 4143: 1025, 4144: 1026, 4145: 1027, 4146: 1028, 4147: 1029, 4148: 1030, 4149: 1031, 4150: 1032, 4151: 1033, 4152: 1034, 4153: 1035, 4154: 1036, 4155: 1037, 4156: 1038, 4157: 1039, 4158: 1040, 4159: 1041, 4160: 1042, 4161
```

(If the code is not clearly visible here , you can find the same code in the github repo, Link for the git hub repo is also provided).

```
new_data = pd.DataFrame({
    'type': ['house'],
    'state': ['ca'],
    'baths': [1],
    'beds': [2],
    'sqfeet': [1000]
})
new_data_cat = ohe.transform(new_data[['type', 'state']])
new_data_cat = pd.DataFrame(new_data_cat.toarray(), columns=ohe.get_feature_names_out(['type', 'state']))
new_data_num = new_data[['baths', 'beds', 'sqfeet']]
new_data_num = scaler.transform(new_data_num)
new_data_num = pd.DataFrame(new_data_num, columns=['baths', 'beds', 'sqfeet'])
new_data = pd.concat([new_data_cat, new_data_num], axis=1)
prediction = rf.predict(new_data)
print('Predicted rent:', prediction)
```

The output files stored are subjected to Bigquery using Amazon Athena to display the tables depending upon the desired output.





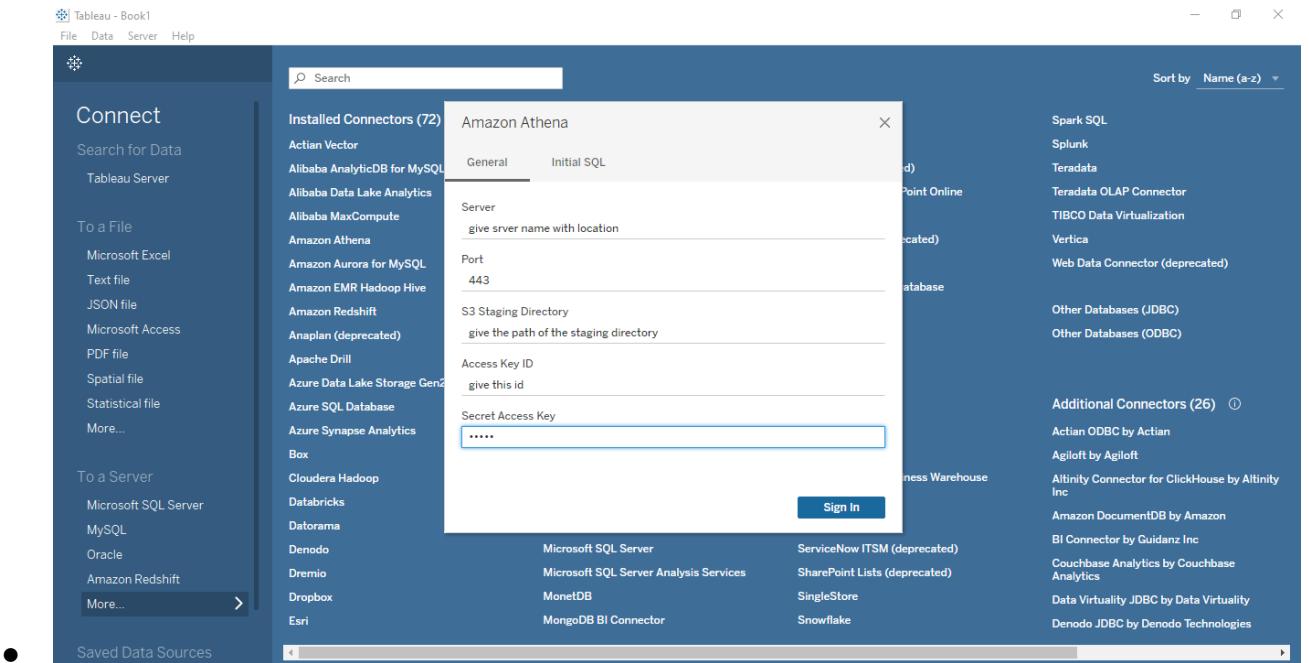
Results (100)

Copy Download results

Search rows

s	parking_options	state	price	cats_allowed	dogs_allowed
south jersey	https://southjersey.craigslist.org		1080	856	
south jersey	https://southjersey.craigslist.org		1445	745	
south jersey	https://southjersey.craigslist.org		980	850	
south jersey	https://southjersey.craigslist.org		1430	750	
south jersey	https://southjersey.craigslist.org		1305	899	
south jersey	https://southjersey.craigslist.org		1495	847	
south jersey	https://southjersey.craigslist.org		1495	847	
south jersey	https://southjersey.craigslist.org		1495	847	
south jersey	https://southjersey.craigslist.org		1515	1150	

- The Athena is then connected to Tableau, a visualization tool to generate charts or reports depending upon the tables generated in the Athena.



The process of running this project is briefly explained in Steps to Run the Application.

TABLEAU LINK for public access :

<https://public.tableau.com/app/profile/manikrishna.sanganabatla>

INSTRUCTIONS FOR DEPLOYMENT:

1. Create an EMR Cluster:

- Log in to the AWS Management Console and navigate to the Amazon EMR service.
- Click on "Create cluster" and configure the cluster settings according to the requirements. Make sure to select the appropriate version of EMR and Spark.
- Specify the number of instances, instance types, and other cluster details. You can also customize advanced options if necessary.
- Once the configuration is complete, click on "Create cluster" to launch the EMR cluster.

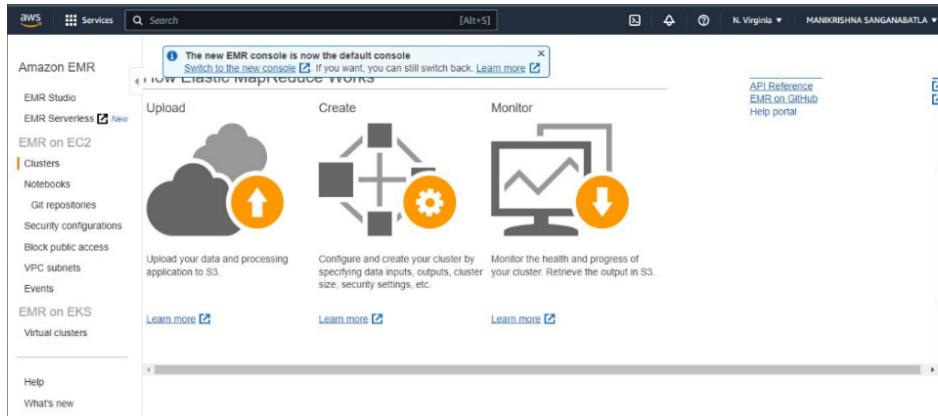
2. Upload Input Data to S3:

- Log in to the AWS Management Console and navigate to the Amazon S3 service.
- Create a new S3 bucket or select an existing one to store the input data.

- Upload the input data files (e.g., CSV, JSON) to the S3 bucket you created. Take note of the S3 bucket path for future reference.
3. Prepare the Spark Job:
- Package the Spark job code and dependencies into a Python file. Make sure all the necessary dependencies are included.
 - If you're using a Python script, ensure it has the necessary import statements and handles the input and output paths correctly.
 - Test the Spark job code locally to ensure it's working as expected before deploying it to the EMR cluster.
4. Connect to the EMR Cluster Master Node:
- Connect to the EMR cluster's master node using SSH. You can use the EC2 key pair associated with the cluster.
 - Once connected, navigate to the Spark directory where the application code is stored.
 - Use this command :
`Ssh -i <location (or) path of the aws key pair> <cluster provided command for integrating master node with local >`
5. Run the Spark Job:
- Submit the Spark job using the spark-submit command, specifying the path to the application code and any necessary arguments. For example:
 - Bash
 - Spark-submit <filename.py>
 - Monitor the job execution logs to ensure it completes successfully.
6. Store Output Files in S3:
- After the Spark job completes, verify that the output files are generated and available in the specified S3 output path. You can check the S3 bucket you specified for the output files.
- (6.1). ML prediction : For this project random forest works better, goto the jupyter notebook or google collab and execute the ml algorithm to get mse, accuracy, give four unique attributes to predict the rent of the house types in various regions. (shift + enter to excute the code).
7. Set Up Athena and Create a Table:
- Go to the AWS Management Console and navigate to the Amazon Athena service.
 - Click on "Query Editor" to open the Athena query editor.
 - Create a new database if required, and then create a table in Athena that matches the schema of the output data.

- Use the "Create Table" wizard or write SQL statements to define the table structure and specify the S3 location of the output data files.
8. Write SQL Queries for Visualizations in Tableau:
- Open Tableau and establish a connection to Athena.
 - Create a new worksheet or dashboard and drag the required dimensions and measures onto the canvas.
 - Write SQL queries using Tableau's interface to retrieve the data from Athena for the visualizations.
 - Customize the visualizations, apply filters, and create interactive elements as needed.
 - Save and publish the visualizations to share with others.

STEPS TO RUN THE APPLICATION:



1. Create an EMR Cluster :

- Sign in to the AWS Management Console and navigate to the Amazon EMR service.
- Click on "Create cluster" and configure the cluster settings according to the requirements. Select the appropriate version of EMR and Spark.
- Specify the number of instances, instance types, and other cluster details. You can also customize advanced options if necessary.
- Click on "Create cluster" to launch the EMR cluster.

2.Upload Data to S3:

- Sign in to the AWS Management Console and go to the Amazon S3 service.
- Create a new S3 bucket or choose an existing one to store the input data.
- Upload the input data files (e.g., CSV, JSON) to the S3 bucket you created. Remember the S3 bucket path for future reference.

3. Connect to the EMR cluster master node :

- Establish an SSH connection to the master node of the EMR cluster using the EC2 key pair associated with the cluster.

4.Transfer Spark Job Code to the Master Node:

- Copy the Spark job code (e.g., my_spark_job.py) to the master node using the SCP or SFTP command. Alternatively, you can clone a Git repository containing the Spark job code directly on the master node.

5.RUN THE SPARK JOB :

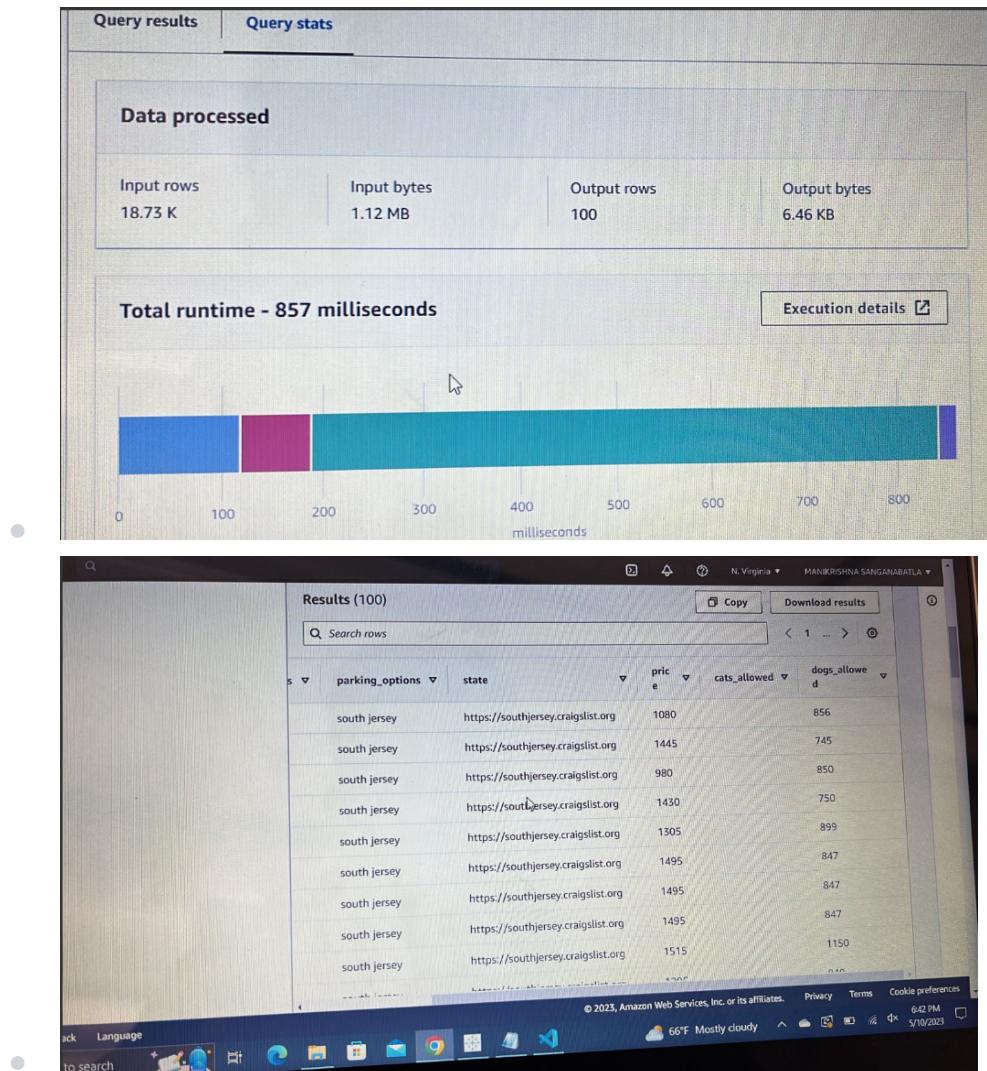
- Execute the Spark job on the master node using the spark-submit command. Adjust the command according to the specific job requirements:
- Copy code
- `Ssh -i <location (or) path of the aws key pair> <cluster provided command for integrating master node with local >`
- `Spark-submit <filename.py>`
- Monitor the job execution logs to ensure it runs successfully.

6.Store Output Files in S3:

- Once the Spark job completes, verify that the output files are generated.
- Upload the output files to the designated S3 output path or create a new S3 bucket for the output if needed.
- (6.1). ML prediction : For this project random forest works better, goto the jupyter notebook or google collab and execute the ml algorithm to get mse, accuracy, give four unique attributes to predict the rent of the house types in various regions. (shift + enter to excute the code).

7.Set Up Athena and Create a Table:

- Access the AWS Management Console and navigate to the Amazon Athena service.
- Open the Athena query editor.
- Create a new database if necessary and define a table that matches the schema of the output data.
- Use the "Create Table" wizard or write SQL statements to specify the table structure and provide the S3 location of the output data files.



8. Write SQL Queries for Visualizations in Tableau:

- Launch Tableau and establish a connection to Athena.
- Create a new workbook and add a new worksheet or dashboard.
- Drag and drop the desired dimensions and measures onto the canvas.
- Compose SQL queries within Tableau to retrieve the data from Athena for the visualizations.

- Customize the visualizations, apply filters, and create interactive elements to suit the analysis requirements.

TEST RESULTS :

1.EMR Cluster Creation:

- Verify that the EMR cluster is successfully created without any errors or failures.
- Check the EMR cluster's status and ensure it transitions to the "Waiting" state.

2.Uploading Input Data to S3:

- Confirm that the input data files are uploaded to the specified S3 bucket.
- Validate the data files' integrity and ensure they are accessible.

3.Running the Spark Job:

- Connect to the EMR cluster's master node using SSH.
- Submit the Spark job and monitor the logs for any exceptions or errors.
- Check the output files generated by the Spark job and confirm they are saved to the specified S3 output path.
- Also check with the “Application User Interfaces”, such as , Spark History server,Yarn Timeline server:

Spark History Server							
Event log directory: s3a://prod.us-west-1.appinfo/src/j-2Y1L1U5I1SVQV/sparklogs							
Last updated: 2023-05-14 22:52:22							
Client local time zone: America/Los_Angeles							
App ID	App Name	Started	Completed	Duration	Spark User	Last Updated	Event Log
application_1683934462317_0003	ADBMSProject	2023-05-12 17:00:32	2023-05-12 17:01:32	1.0 min	hadoop	2023-05-12 17:02:47	<button>Download</button>
application_1683934462317_0002	ADBMSProject	2023-05-12 16:43:00	2023-05-12 16:43:46	47 s	hadoop	2023-05-12 16:44:47	<button>Download</button>
application_1683934462317_0001	ADBMSProject	2023-05-12 16:39:55	2023-05-12 16:40:43	48 s	hadoop	2023-05-12 16:40:47	<button>Download</button>

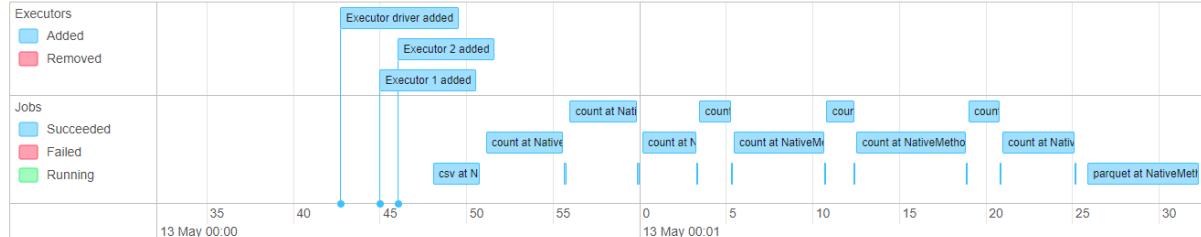
Showing 1 to 3 of 3 entries
Show incomplete applications

Spark Jobs (?)

User: hadoop
 Total Uptime: 1.0 min
 Scheduling Mode: FIFO
 Completed Jobs: 20

Event Timeline

Enable zooming



Completed Jobs (20)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
19	parquet at NativeMethodAccessoimpl.java:0 parquet at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:25	6 s	1/1	8/8
18	count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:25	26 ms	1/1 (1 skipped)	1/1 (8 skipped)

Completed Jobs (20)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
19	parquet at NativeMethodAccessoimpl.java:0 parquet at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:25	6 s	1/1	8/8
18	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:25	26 ms	1/1 (1 skipped)	1/1 (8 skipped)
17	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:20	4 s	1/1	8/8
16	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:20	25 ms	1/1 (1 skipped)	1/1 (8 skipped)
15	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:18	2 s	1/1	8/8
14	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:18	30 ms	1/1 (1 skipped)	1/1 (8 skipped)
13	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:12	6 s	1/1	8/8
12	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:12	26 ms	1/1 (1 skipped)	1/1 (8 skipped)
11	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:10	2 s	1/1	8/8
10	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:10	27 ms	1/1 (1 skipped)	1/1 (8 skipped)
9	count at NativeMethodAccessoimpl.java:0 count at NativeMethodAccessoimpl.java:0	2023/05/13 00:01:05	5 s	1/1	8/8
-	-	-	-	-	-

4. Storing Output Files in S3:

- Verify that the output files are present in the designated S3 output path.
- Ensure the output files contain the expected data and have the correct file format.
- ML prediction results are displayed as below :

```

rf.fit(X_train, y_train)
# Make predictions on the test set
rf_pred = rf.predict(X_test)

# Evaluate the model using Mean Squared Error
rf_mse = mean_squared_error(y_test, rf_pred)
print('Random Forest MSE:', rf_mse)

# Make predictions on new data
new_data = pd.DataFrame({
    'type': ['house'],
    'state': ['ca'],
    'baths': [1],
    'beds': [2],
    'sqfeet': [1000]
})
new_data_cat = ohe.transform(new_data[['type', 'state']])
new_data_cat = pd.DataFrame(new_data_cat.toarray(), columns=ohe.get_feature_names_out(['type', 'state']))
new_data_num = new_data[['baths', 'beds', 'sqfeet']]
new_data_num = scale.transform(new_data_num)
new_data_num = pd.DataFrame(new_data_num, columns=['baths', 'beds', 'sqfeet'])
new_data = pd.concat([new_data_cat, new_data_num], axis=1)
prediction = rf.predict(new_data)
print('Predicted rent:', prediction)

```

Fitting 5 Folds for each of 17 candidates, totalling 135 fits
Best Parameters: {'max_depth': 5, 'min_samples_split': 2, 'n_estimators': 200}

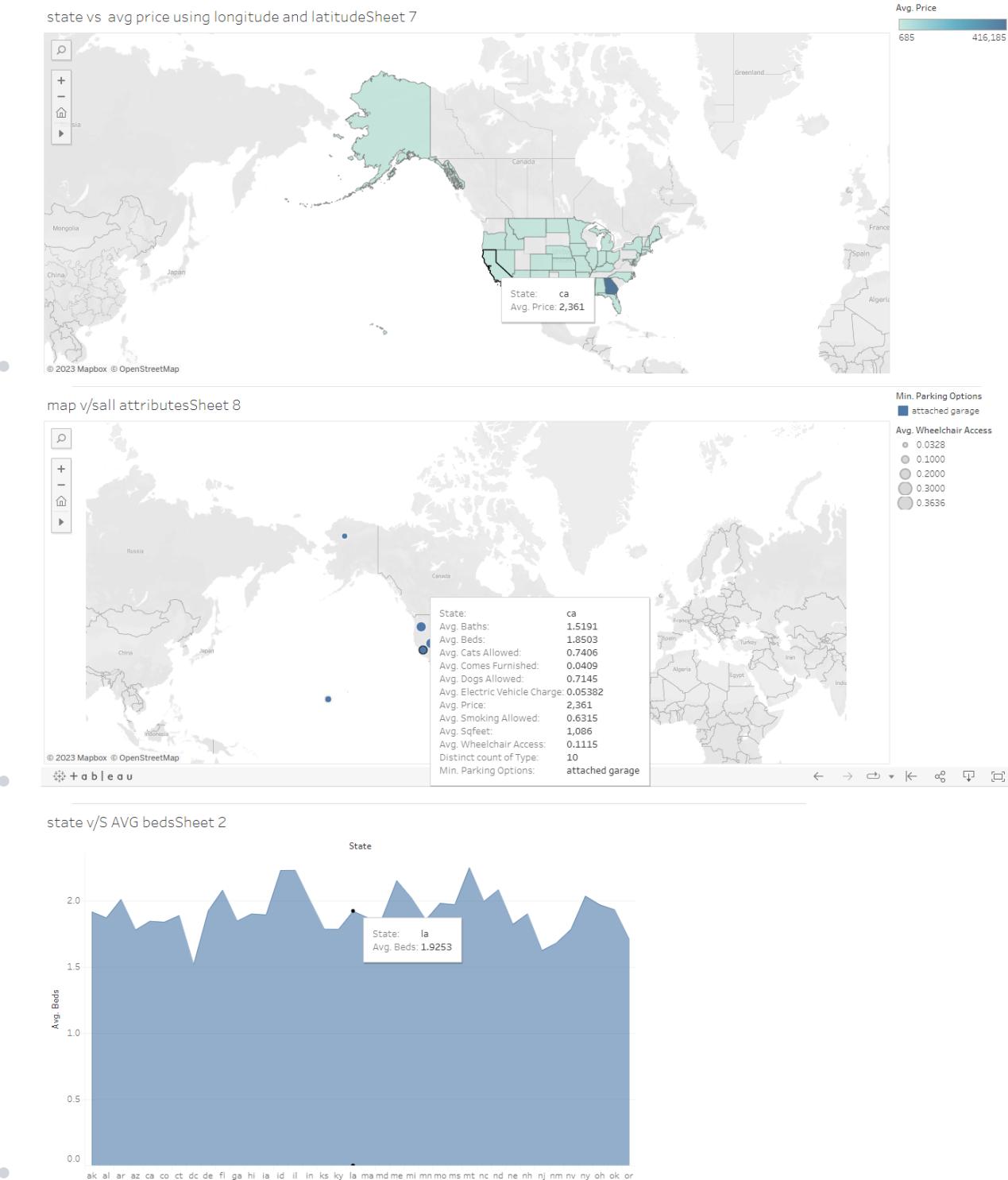
Predicted rent: [1929.28699387]

5. Using Output Data in Athena:

- Create a table in Athena using the output data files stored in S3.
- Run sample SQL queries against the table to retrieve and validate the data.
- Ensure the data matches the expected format and aligns with the original input.

6. Writing SQL Queries for Visualizations in Tableau:

- Connect Tableau to Athena and establish a successful connection.
- Create visualizations using the SQL queries against the Athena table.
- Verify that the visualizations accurately represent the analyzed data and provide the desired insights.



These are general steps to consider for testing the house rental analysis project. It's important to customize and expand upon these tests based on the specific requirements and project functionalities. Additionally, conducting end-to-end integration tests that cover the entire

workflow, from data ingestion to visualization, is recommended to ensure seamless execution and accurate results.

Note : Major part of the above execution of this project can be found in the Zip file which is uploaded in the canvas.

REFERENCES :

<https://www.kaggle.com/datasets/rkb0023/houserentpredictiondataset>

(All the information mentioned in the last deployment and steps to run methods also can be found in AWS Educate website)

<https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-gs.html>

(same above doc for S3 bucket creation, Athena)

https://help.tableau.com/current/pro/desktop/en-us/examples_amazonathena.htm