

REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM

21CSE533J - Advanced - Data Processing Techniques Mini Project Report

Submitted by

JUJJAVARAPU SUJAN CHOWDARY (RA2412033010001)

GANDAM SAI RAM (RA2412033010034)

DILEEP SAI ANDE (RA2412033010026)

NAVEEN KUMAR REDDY (RA2412033010029)

Under the Guidance of

Dr. S Prabakeran

Associate Professor, Department of Networking and Communications

In Partial fulfillment of the requirements for the degree of

**MASTER OF TECHNOLOGY
in
CLOUD COMPUTING**



**DEPARTMENT OF NETWORKING AND COMMUNICATIONS
COLLEGE OF ENGINEERING AND TECHNOLOGY
SRM INSTITUTE OF SCIENCE AND TECHNOLOGY
(Under Section 3 of UGC Act, 1956)
SRM NAGAR, KATTANKULATHUR - 603 203
CHENGALPATTU DISTRICT
MAY 2025**

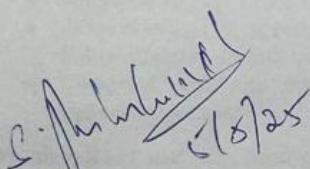


SRM INSTITUTE SCIENCE AND TECHNOLOGY

KATTANKULATHUR – 603 203

BONAFIDE CERTIFICATE

Certified that this M.Tech, Minor Project report titled "**REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM**" is the bona fide work of **Jujjavarapu Sujan Chowdary (RA2412033010001)**, **Gandam Sai Ram (RA2412033010034)**, **Dileep Sai Ande (RA2412033010026)**, and **Naveen Kumar Reddy (RA2412033010029)** who carried out the project work under my supervision. Certified further, that to the best of my knowledge the work reported here does not form part of any other thesis or dissertation based on which a degree or award was conferred on an earlier occasion for this or any other candidate.


Dr. S Prabakeran

SUPERVISOR

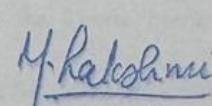
Associate Professor

Department of Networking

and Communications,

SRM INSTITUTE OF

SCIENCE AND TECHNOLOGY, KTR


Dr. M. Lakshmi

PROFESSOR AND HEAD

Department of Networking

and Communications,

SRM INSTITUTE OF

SCIENCE AND

TECHNOLOGY, KTR



SRM

INSTITUTE OF SCIENCE & TECHNOLOGY
Decreed to be University u/s 3 of UGC Act, 1956

DEPARTMENT OF NETWORKING AND COMMUNICATIONS

SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

OWN WORK DECLARATION FORM

DEGREE COURSE: M-TECH - CLOUD COMPUTING

STUDENT NAMES: JUJJAVARAPU SUJAN CHOWDARY, GANDAM SAI RAM, DILEEP SAI ANDE,
NAVEEN KUMAR REDDY

REGISTRATION NO.: RA2412033010001, RA2412033010034, RA2412033010026, RA2412033010029

TITLE OF WORK: - REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM

We hereby certify that this assessment complies with the University's Rules and Regulations relating to Academic misconduct and plagiarism, as listed in the University Website, Regulations, and the Education Committee guidelines. We confirm that all the work contained in this assessment is my / our own except where indicated, and that We have met the following conditions:

- Clearly referenced / listed all sources as appropriate
- Referenced and put in inverted commas all quoted text (from books, web, etc)
- Given the sources of all pictures, data etc. that are not my own
- Not made any use of the report(s) or essay(s) of any other student(s) either past or present
- Acknowledged in appropriate places any help that I have received from others (e.g. fellow students, technicians, statisticians, external sources)
- Compiled with any other plagiarism criteria specified in the Course handbook / University website

We understand that any false claim for this work will be penalized in accordance with the university policies and regulations.

DECLARATION:

We are aware of and understand the University's policy on Academic misconduct and plagiarism and I certify that this assessment is our own work, except where indicated by referring, and that we have followed the good academic practices noted above.

JuJJavarapu Sujan Chowdary (RA2412033010001)

Gandam Sai Ram (RA2412033010034)

Dileep Sai Ande (RA2412033010026)

Naveen Kumar Reddy (RA2412033010029)

If you are working in a group, please write your registration numbers and sign the date for every student in your group.

ACKNOWLEDGEMENT

We express our humble gratitude to Dr. C. Muthamizhchelvan, Vice-Chancellor, SRM Institute of Science and Technology, for the facilities extended for the project work and his continued support.

We extend our sincere thanks to Dean-CET, SRM Institute of Science and Technology, Dr.T.V. Gopal, for his invaluable support.

We wish to thank Dr. Revathi Venkataraman, Professor & Chairperson, School of Computing, SRM Institute of Science and Technology, for her support throughout the project work.

We are incredibly grateful to our Head of the Department, Dr. M . Lakshmi , Professor and Head, Department of Networking and Communications, School of Computing, SRM Institute of Science and Technology, for her suggestions and encouragement at all the stages of the project work.

Our inexpressible respect and thanks to our guide, Dr. S. Prabakeran, Associate Professor, Department of Networking and Communications, SRM Institute of Science and Technology, for providing us with an opportunity to pursue our project under his mentorship. He provided us with the freedom and support to explore the research topics of our interest. His passion for solving problems and making a difference in the world has always been inspiring.

We sincerely thank the Networking and Communications, Department staff and students, SRM Institute of Science and Technology, for their help during our project. Finally, we would like to thank parents, family members, and friends for their unconditional love, constant support, and encouragement.

Jujjavarapu Sujan Chowdary [RA2412033010001]

Gandam Sai Ram [RA2412033010034]

Dileep Sai Ande [RA2412033010026]

Naveen Kumar Reddy [RA2412033010029]

ABSTRACT

This project presents an AI-powered real estate app to improve property discovery, price forecasting, and purchasing or leasing choices for buyers and lessees. The conventional real estate websites lack intelligent search results, correct price data, and customized recommendations, and thus result in suboptimal property choice and restricted market intelligence. This project overcomes these shortcomings by developing a high-performance, cloud-native solution on the basis of an ensemble of sophisticated machine learning algorithms, big data processing, and generative AI.

At the center, the project utilizes Google Cloud Vertex AI to train machine learning models to forecast real estate prices and rental prices based on various property features such as area, location, amenities, and property type. PySpark is employed to set up a data pipeline in AWS SageMaker Studio for enterprise-level data transformation, feature engineering, and enrichment of property data. Amazon S3 stores and manages data with durability and scalability.

Processed data is presented graphically through Amazon QuickSight, with individual dashboards for buying and leasing real estate. The dashboards include calculated fields such as price per square foot, price difference between forecast and actual, and easy-to-use filters. Privacy is maintained by removing owner-sensitive data.

For enabling the system to be intelligent and interactive, the Gemini API and ChatGPT API are implemented in such a manner that the users are allowed to interact with the system as natural language inquiries. This provides intuitive searching as well as real-time insights and provides end-users with a conversational interface.

The result is a secure, scalable, and smart real estate app that gives users precise pricing predictions, data-driven recommendations, and a simple-to-use interface. It shows the power of leveraging the newest cloud platforms and AI capabilities to upend conventional businesses such as real estate.

TABLE OF CONTENTS

	Page.No
ABSTRACT.....	
LIST OF TABLES.....	
LIST OF FIGURES.....	
CHAPTER : 1 Introduction.....	1
1.1 Motivation.....	1
1.2 Problem Statement.....	1
1.3 Objective of the project.....	2
1.4 Scope.....	2
1.5 Relevance.....	3
CHAPTER : 2 Literature Review.....	4-6
CHAPTER : 3 Proposed Methodology.....	7
3.1 Project Approach.....	7
3.2 Data Collection.....	7
3.2.1 DataSet Description.....	8
3.3 Infrastructure Setup and Tool Selection.....	10
3.3.1 AWS Infrastructure Setup.....	11
3.3.2 Tool Selection and Configuration.....	11
3.4 Machine Learning Model Development.....	12
3.5 Customer Page.....	19
3.5.1 Gemini Setup.....	20
3.5.2 AWS S3 Setup.....	21
3.5.3 Get_Command_input() Function.....	21
3.5.4 Fetch_properties() Function.....	21
3.5.5 Get_Gemini_Insight() Function.....	22
3.5.6 Main() Function-The APP entry point.....	23
3.6 Owner.....	24
3.6.1 Libraries and Initialization.....	25
3.6.2 Model Loading.....	26
3.6.3 User interface(UI).....	26
3.6.4 From Input Collection.....	26

3.6.5 Prediction via SparkML Models.....	27
3.6.6 Open AI Chat Insight...	27
3.6.7 Sample Output.....	28
CHAPTER : 4 Implementation	30
4.1 System Design.....	30
4.2 Coding and Tools.....	30
4.3 Process Flow...	31
4.4 Challenges and obstacles.....	32
4.5 Cloud Development and Integration(AWS + GCP).....	32
4.5.1 AWS Components.....	32
4.5.1.1 AWS S3-Centralized Data Storage.....	32
4.5.1.2 AWS Sagemaker studio-Model Visualization & Experimentation.....	33
4.5.1.3 AWS QuickSite-Business Intelligence Dashboard.....	34
4.5.2 GCP Components.....	36
4.5.2.1 Vertex AI - Model Development and Deployment.....	36
4.5.2.2 Vertex Pipeline.....	36
4.5.3 Model Cloud Architecture: WHY & HOW.....	37
4.5.4 Multi Cloud Strategies Highlights.....	37
4.5.5 Security and Access Control.....	37
CHAPTER : 5 Results and Discussions.....	39
5.1. Customer Dashboard Insights.....	39
5.2. Owner Dashboard Insights.....	42
5.3 Visual Insights with AWS Quick sight.....	44
CHAPTER : 6 Conclusion	48
CHAPTER : 7 FUTURE WORK	49
References.....	50
Project Files and Demo References	52
Appendix.....	53

LIST OF TABLES

TABLE NO.

PAGE.NO

Table 2.1	Cloud services and functionalities	4
Table 2.2	Integration of LLM and cloud components	5
Table 3.1	Dataset Schema	10
Table 3.2	Infrastructure Setup Services	11
Table 3.3	Multi-Cloud Tools	11
Table 3.4	Model Evaluation Metrics	19
Table 4.1	Challenges and Obstacles	32
Table 4.2	Detailing of Components	37
Table 5.1	Filtered Results	40

TABLE OF FIGURES

FIGURE.NO		PAGE NO
Figure 3.1	Architecture Diagram	8
Figure 3.2	Loading Dataset	12
Figure 3.3	Initialize Spark Session	13
Figure 3.4	Lowering & Trimming	13
Figure 3.5	Replacing Strings	14
Figure 3.6	Filling Missing Values	14
Figure 3.7	Splitting Dataset	15
Figure 3.8	Feature Engineering	15
Figure 3.9	Feature columns	16
Figure 3.10	ML Pipeline	17
Figure 3.11	Testing and Training	17
Figure 3.12	Model Evaluation Metrics	18
Figure 3.13	Price Prediction	18
Figure 3.14	Customer flow chart	20
Figure 3.15	Gemini API key	20
Figure 3.16	AWS S3 setup	21
Figure 3.17	Fetching Properties	22
Figure 3.18	Plain Text Report	23
Figure 3.19	Main Function Setup	24
Figure 3.20	Owner Dashboard	25
Figure 3.21	Model Loading	26

Figure 3.22	Prediction via Spark ML	27
Figure 3.23	Open AI Chat Insight	28
Figure 3.24	Sample Output	29
Figure 4.1	AWS S3 Setup	33
Figure 4.2	SageMaker Studio Setup	33
Figure 4.3	Jupyter Lab in Sage maker Studio	34
Figure 4.4	SageMaker JupyterLab HomePage	34
Figure 4.5	Quicksite Buyer Dashboard	35
Figure 4.6	Quicksite Rental Dashboard	35
Figure 4.7	Google Vertex AI	36
Figure 4.8	Vertex AI Homepage	37
Figure 4.9	Login Page	38
Figure 5.1	Customer Result page	41
Figure 5.2	Customer LLm recommend response	42
Figure 5.3	Owners model house prediction and LLm recommended response	44
Figure 5.4	Average Rental Month Cost by Locality & House Type	45
Figure 5.5	Average Rent Month Cost by Locality & BHK	45
Figure 5.6	Total Rent Price per Year by Locality	46
Figure 5.7	Average Rent Cost by Availability Date,BHK & Locality	46
Figure 5.8	Average Rent Month Cost by Availability Date, Locality & House type	47
Figure 5.9	Average Rent Month Cost by House Type & Locality	47

CHAPTER-1

INTRODUCTION TO REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM

Real estate businesses produce huge daily volumes of operational data from listings and market trends along with buyer demands and legal files. Real estate data volumes persistently go unused due to systems that cannot perform analytical functions nor deliver personalized outputs. Intelligent systems need to address the current real estate data processing issues as digital adoption gains momentum along with changing consumer needs. Scalable intelligent real estate solutions based on cloud computing technologies and machine learning and large language model frameworks deliver useful insights and increase user interaction and quick data-driven decisions.

1.1 Motivation:

Modern users find that traditional real estate platforms fall short because they present minimal analytical tools as well as a deficient level of interactive performance that satisfies contemporary user needs. Property decision-making becomes more challenging for buyers, renters and agents because of market instability together with price volatility and unorganized property data. The project emerges from the necessity to handle crucial market barriers through multi-cloud architecture integrated with generative AI functionalities. Through connected data pipelines and predictive modeling and conversational AI integration the system creates a new way for users to interact with property data. The objective targets developing a system that is technically superior while remaining centered on users' needs and adapts to market changes.

1.2 Problem Statement:

The present-day real estate market creates multiple obstacles that influence home purchasers and lessees at all periods of their property search and evaluation process. Existing real estate platforms achieve minimal success because they show standard properties without AI analytics or individual property recommendation capabilities. Traditional real estate methods conduct temporary price assessments from past sale data that leads to bad analyses along with incorrect

evaluations and higher operational costs. Agents working in the real estate sector today face a challenge because their industry needs AI capabilities to analyze market trends and determine property values which adjusts customer outreach.

1.3 Objective of the Project:

This project aims to construct an AI-based real estate platform that raises user satisfaction through its intelligent search function added to its price and rental value estimation and interactive query system. The platform is designed to:

1. The application enables ML models to predict real estate property prices and rental value assessments based on training datasets.
2. Distributed processing through PySpark enables the transformation and engineering of scalable data
3. The user should access clean structured data visualizations which appear through intuitive dashboard interfaces.
4. The system uses integrated generative AI APIs from Gemini and ChatGPT so users can conduct natural language interactions.
5. Decision-making abilities of users should be enhanced through the system because they search for properties or set prices and analyze market trends.

1.4 Scope:

This project entails full development of a real estate analytics system that incorporates cloud computing and artificial intelligence functionalities. The scope includes:

- The implementation of AWS uses PySpark deployed through SageMaker Studio to handle property data acquisition and transformation involving buy and rent operations.
- The implementation of ML models to forecast prices and rents took place through the deployment of models on Google Cloud Vertex AI.
- The solution stores all intermediate data and final data products in the Amazon S3 data storage system.
- A set of Amazon QuickSight dashboards assists visualization of predictions by combining metrics and trends.

- The system uses Gemini and ChatGPT APIs to enable users to interact conversationally with the platform when searching for properties and receiving recommendations. The solution covers metropolitan areas in India and contains standard properties that combine area and BHK type information with location and amenities details and property classification.

1.5 Relevance:

Decisions about real estate matter involve handling both social risks together with lengthy financial commitments. Processing accurate price data along with market analytics and specific recommendations guides buyers and sellers and tenants toward improved confidence about their decisions. Users at present can access systems that give them friendly interfaces to understand alongside non-expert users through generative AI programs. Technical excellence of this project stands out because machine learning systems operate on real estate data alongside cloud scalability to deliver data-driven personalized interactions within transparent operations.

CHAPTER 2

LITERATURE REVIEW FOR REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM

Real estate analytics underwent substantial transformation when data engineering united with cloud computing and large language models (LLMs). Real estate dataset unifies various heterogeneous large scales of information from documentation records to location information together with imaging and textual documents from listings and legal files. AI-enabled multi-cloud platforms function as a system to establish effective data pipeline architectures for complicated scalable applications [1] [3].

Infrastructure services accessible through GCP and AWS platforms simplify these operations. Vertex AI at Google Cloud Platform delivers an integrated platform that enables users to handle model development training together with deployment services through its machine learning (ML) model building service. AutoML along with pipeline management enables the software to quicken the development process and boost operational effectiveness for real estate forecasting and classification activities [2] [4]. AWS functions as an outstanding base for data handling and storage services. Amazon S3 combines trusted data storage features that enable usage with AWS SageMaker and its Apache Spark component enables quick processing of extensive datasets [5] [6]. The processed data enables users to build business intelligence by developing self-service dashboards through AWS QuickSight [6, 9].

Cloud Provider	Service	Functionality
GCP	Vertex AI	ML model development and deployment
AWS	S3, SageMaker	Data storage and scalable Spark-based transformation
AWS	QuickSight	Visualization and interactive dashboards

Table:2.1 Cloud services and functionalities

Real estate data processing has experienced major advancements with the creation of GPT from OpenAI and Gemini from Google which are large language models. These NLP models stand as leading information revelation systems which extract regulated data from listings and contracts while demonstrating intellectual understanding and document abstraction features [10] [11] [13]. GPT functions as a tool to help agents along with investors to speed up their choices by identifying pricing patterns and locating legal requirements within processed contracts [14].

The development frameworks of Streamlit enable programmers to embed LLMs in their applications for interface design purposes. Streamlit enables developers to create basic web-based user interfaces that display immediate text-based and analytical visualizations of predictions together with text analysis outputs through interactive dashboards [15, 16]. The built-in interfaces in this display system allow workers with non-technical skills to understand analytical information [17].

Research evidence indicates that real estate domains gain maximum value from the combination of cloud-native platforms with interactive visual tools and LLM-based intelligence [18, 19, 20]. The integration of different cloud technologies and language models enables endless future analysis while enabling the detection of fraud activities.

The combined delivery method transforms to changing business conditions by employing intelligent data processing in conjunction with automated workflows. A multi-cloud solution enables the system to handle operational effectiveness features along with redundancy while adhering to cost requirements [8]. User requests enter deep conceptual processing via LLM capabilities that create automated reports as well as control client interactions through chatbots for improved digital interactions with end-users [13, 16]. Scalable cloud infrastructure and the combination of LLMs and ML enables developers to create sustainable systems for real estate property data management that facilitate smart city development and IoT monitoring systems [20].

Various investigations confirm how AI ethical guidelines and cross-cloud connection methods offer different positive and negative aspects. The implementation of API layer combination with governance protocols resolves latency and compliance issues between GCP and AWS ecosystems according to [3] and [7]. Federated learning acts as a solution that creates distributed cloud model training frameworks which preserve data privacy [12].

Challenge	Technology Used	Benefit
Unstructured text (e.g., listings, legal)	LLMs (GPT, Gemini)	Semantic parsing and entity recognition
Data volume and variety	AWS S3, SageMaker	Scalable ETL processing
Visualization and insights	QuickSight, Streamlit	Dynamic dashboards, interactive reports
Real-time model inference	Vertex AI	Fast deployment of ML pipelines

Table:2.2 Integration of LLM and cloud components

LLMs help predictive analytics when used for building AVMs that support financial underwriting operations and risk management systems [19]. Use of ML models that bridge transactional and demographic data alongside systems produces accuracy at higher levels than current traditional practices [11].

Real estate analytics will see revolutionary transformations because multi-cloud services unite with LLMs according to [7]. A text analysis system with LLM functions operates through a modular pipeline that has been developed by joint GCP model development and AWS scalable visualization and transformation capabilities [12]. The infrastructure enhances operations by providing data-driven smart insights while establishing better stakeholder relations from property development to management phases [14].

CHAPTER - 3

PROPOSED METHODOLOGY OF REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM

The system establishes a real estate price prediction solution through the integration of Google Cloud and AWS services. The system reaches top performance and flexible scalability because it merges both Google Cloud and AWS management services. The complete method begins with data retrieval followed by data preparation after which a machine learning model is developed before the predictions stage before generating visual outputs. Each component inside the system structure benefits from a modular design which enables independent management until the point of individual update capabilities. The system delivers accurate estimates of rent and purchase prices while allowing users to enhance decision quality through its interactive interface. The system operates with updated tools that incorporate Vertex AI together with AWS S3 and QuickSight for data storage and visualization and Gemini and ChatGPT LLMs for automatic code generation and advanced querying functionalities.

3.1 Project Approach

The project implements a hybrid cloud model which chooses optimal features between AWS and GCP. The housing data first rests in the AWS S3 storage system because it provides scalable durable and accessible storage. The data preprocessing and feature engineering tasks run through Python execution on AWS SageMaker Studio Lab. The refined data gets transferred from the data refinement process to Vertex AI within Google Cloud for machine learning model training operations from the fig 3.1. The XGBoost and Linear Regression regression models act as predictors to calculate rental and selling prices of housing properties. Real-time or batch prediction occurs through Vertex AI endpoints which deploy the trained models. The final step involves representing the results through AWS QuickSight dashboards.

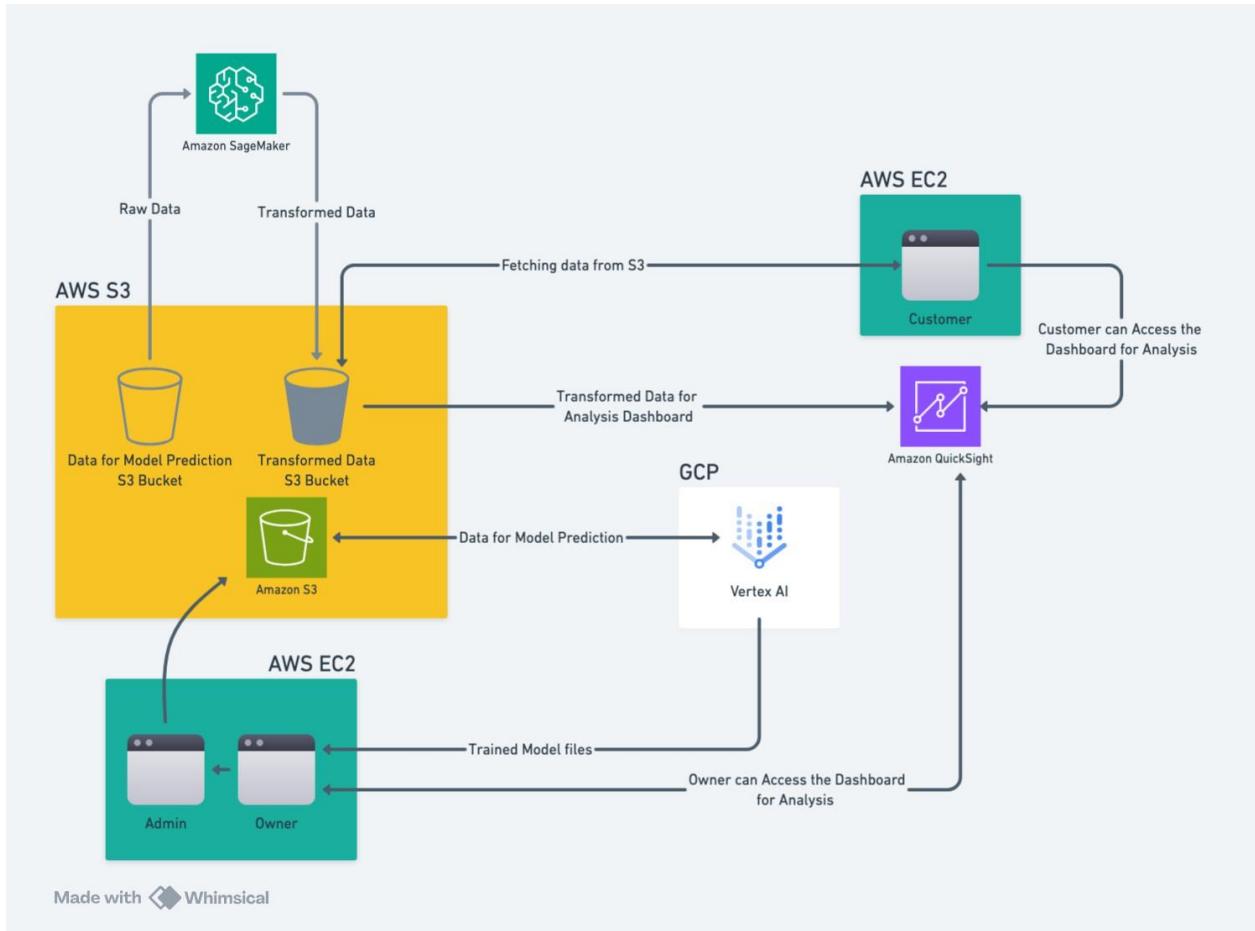


Fig: 3.1 Architecture Diagram

3.2 Data Collection

The core plan for this initiative builds a data-based system that offers predictive insights and recommendation solutions supported by a hybrid cloud solution. Through this strategy we can leverage the most beneficial elements of GCP's ML capabilities in Vertex AI plus the effective data transformation services and visualization tools from AWS (SageMaker, S3, QuickSight) with intelligent interaction capabilities from Gemini and ChatGPT APIs.

3.2.1 Dataset Description

Real estate listings from major Indian cities constitute the dataset which came from open-source property platforms that cover both buy and rent segments. This database contains organized property data with fundamental features regarding residential homes that consist of town and district identification alongside BHK setup sizes together with bathroom numbers and furnishing types and property classifications and selling prices and rent costs with vacancy conditions and amenity specifications as shown in the table 3.1.

Column Name	Data Type	Description
Property_ID	String	Unique identifier for each property (e.g., BU12345678 or RE87654321)
City	String	City where the property is located (e.g., Chennai)
Locality	String	Specific area or neighborhood in the city
Property Purpose	String	Indicates whether the property is for 'buy' or 'rent'
House Type	String	Type of house (e.g., apartment, villa, individual house, gated community)
BHK	Integer	Number of bedrooms
House Floor	Integer	The floor on which the property is located
Building Floor	Integer	Total number of floors in the building
Property Age	Integer	Age of the property in years
Facing	String	Direction the property faces (e.g., East, West, North-East, etc.)
Square Feet	Integer	Built-up area in square feet
rent_monthly_cost	Integer	Monthly rent amount for rental properties (INR)
Owner Price per Sqft	Integer	Owner's listed price per square foot (for buying)
GOV Price per Sqft	Integer	Government guideline price per square foot
Security	String	Security feature available (e.g., CCTV, Security Guards, etc.)
Gym	String	Gym facility availability (Yes or No)
Convention Hall	String	Convention or party hall availability
Parking	String	Parking options (e.g., car, bike, both, no parking)
Water Supply	String	Source of water supply (corporation, borewell, or both)
Bathrooms	Integer	Number of bathrooms
Balcony	Integer	Number of balconies
Availability Date	Date	Date when the property becomes available
Deposit Cost	Integer	Refundable security deposit (for rentals)
Furnishing	String	Furnishing status (unfurnished, semi, full)
HOA Fees (₹)	Float	Homeowners Association fees (Maintenance) in INR
Flooring Type	String	Type of flooring (e.g., Marble, Tiles, Wood, Vinyl, etc.)
Roof Type	String	Type of roof structure (e.g., Concrete, Metal)

Property Style	String	Architectural style of the property (e.g., Modern, Traditional)
View	String	Property view (e.g., Sea View, City View, Garden View, None)
Nearby Schools	Integer	Count of nearby schools
Nearby Hospitals	Integer	Count of nearby hospitals
Nearby Supermarkets	Integer	Count of nearby supermarkets
IT Hub Present	String	Nearby IT company or tech park, if any
Swimming Pool	String	Swimming pool availability (Yes or No)
Lot Shape	String	Shape of the plot (e.g., Rectangular, Irregular)
Sewer System	String	Type of sewer system (Public, Septic)
Address	String	Full address of the property
owner_name	String	Full name of the property owner
owner_email	String	Contact email ID of the owner
Owner Phone	String	Contact number of the owner
Pincode	String	Postal code of the property's location
State	String	State in which the property is located (e.g., Tamil Nadu)
Photos	String	Link to property photo or placeholder image

Table:3.1 Dataset Schema

The dataset enters AWS S3 buckets as CSV files before it gets stored. The clean-up process with null value management alongside categorical field encoding uses PySpark in SageMaker Studio. During feature engineering the team develops extra fields including price per square foot measurements and location group clustering. For optimized modeling and visualization purposes the data gets divided into multiple sections based on city and property status (buy/rent). Vertex AI trains ML models using the enriched data that functions as the base for creating dashboards in QuickSight.

3.3 Infrastructure Setup and Tool Selection

A strong cloud infrastructure using AWS and Google Cloud components was built for smooth data processing combined with model development and visual display. The design system creates effective data communication that connects different cloud components through their individual advantages. The storage preprocessing and dashboard creation tasks operate on AWS infrastructure yet Google Cloud is responsible for ML model development and deployment activities. The chosen tools required integration with scalability features and easy connector support along with full compatibility with all cloud platforms.

3.3.1 AWS Infrastructure Setup

The data management and visualization tasks heavily depend on AWS's functional role. The following services were used as shown in the table below 3.2:

Service	Purpose
Amazon S3	Used to store raw and processed housing datasets in CSV format.
SageMaker Studio Lab	Performed data preprocessing, feature engineering, and Python scripting.
Amazon QuickSight	Built interactive dashboards for rent and buy price comparisons and trends.

Table:3.2 Infrastructure Setup Services

The preprocessing operations through SageMaker notebooks accessed S3-stored data before the information was imported into QuickSight for analytics purposes. The IAM roles operated to facilitate safe data exchanges among different services. The arrangement offers flexible secure scalable data management capabilities.

3.3.2 Tool Selection and Configuration

The tools chosen precisely matched the multi cloud and ML and visualization deliverables of the project. Below is a summary as shown in the table below 3.3:

Tool/Service	Platform	Function
Google Vertex AI	GCP	ML model training, tuning, and deployment.
Amazon S3	AWS	Centralized storage for multi-cloud access.
Amazon QuickSight	AWS	Dashboard for visualizing rental/sale predictions.
Gemini & ChatGPT	Cross-platform	Used for code generation and data analysis support.

Table:3.3 Multi-Cloud Tools

Interoperability was a key requirement for which these tools received configuration. A data pipeline architecture was built to provide secure data transfers between AWS and GCP which created a uniform ML process from start to finish.

3.4 Machine Learning Model Development

This section outlines the **full pipeline for developing two machine learning models** using **PySpark** for:

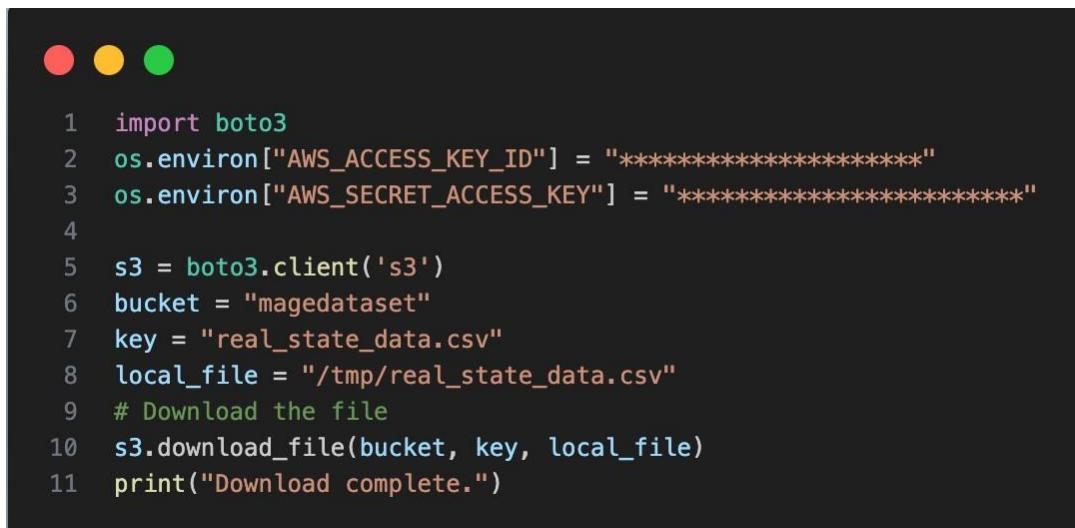
1. **Property Price Prediction** (for buyers)
2. **Monthly Rent Prediction** (for rental properties)

We use **linear regression**, enhanced with **feature engineering** and **categorical encoding**, to model complex patterns in the housing data.

Step 1: Load Dataset from AWS S3

The machine learning workflow starts with accessing the real estate dataset which resides on AWS S3. Secure configuration of AWS credentials enables connection to the specified S3 bucket for downloading the CSV file to a temporary local storage location. The downloaded dataset acts as our base data to support all model development during our project's data processing.

We access and download the dataset from S3 from the below fig 3.2.



```
1 import boto3
2 os.environ["AWS_ACCESS_KEY_ID"] = "*****"
3 os.environ["AWS_SECRET_ACCESS_KEY"] = "*****"
4
5 s3 = boto3.client('s3')
6 bucket = "magedataset"
7 key = "real_state_data.csv"
8 local_file = "/tmp/real_state_data.csv"
9 # Download the file
10 s3.download_file(bucket, key, local_file)
11 print("Download complete.")
```

Fig: 3.2 Loading Dataset

Step 2: Initialize Spark Session and Load CSV

A Spark session launches to allow distributed data processing through the use of PySpark. A Spark DataFrame automatically infers its schema when the downloaded CSV file gets loaded into it. The data exploration of the dataset structure assists us in examining column names together with their corresponding data types so we can perform transformations effectively for modeling purposes from the fig 3.3.



```
● ● ●

1 from pyspark.sql import SparkSession
2 spark = SparkSession.builder.appName("RealEstateData").getOrCreate()
3 df = spark.read.csv(local_file, header=True, inferSchema=True)
4 df.printSchema()
5 df.show(5)
```

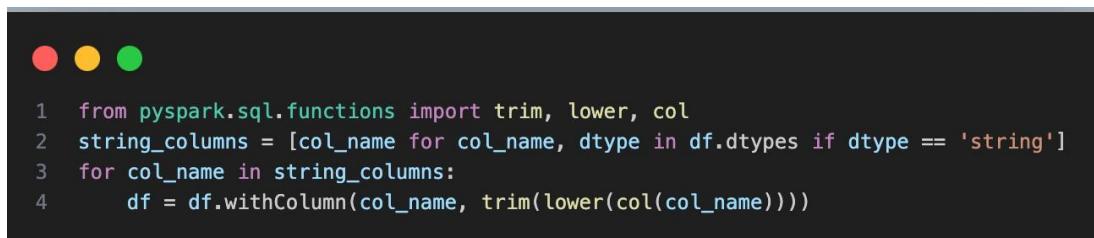
Fig:3.3 Initialize Spark Session

Step 3: Data Cleaning & Preprocessing

To ensure consistent and reliable analysis, the dataset undergoes several cleaning and preprocessing steps:

Lowercase and Trim All String Columns

The normalization technique of text data standardizes field values by making every string lowercase while eliminating extra spaces which reduces data inconsistency that stems from formatting differences as shown in the fig 3.4.



```
● ● ●

1 from pyspark.sql.functions import trim, lower, col
2 string_columns = [col_name for col_name, dtype in df.dtypes if dtype == 'string']
3 for col_name in string_columns:
4     df = df.withColumn(col_name, trim(lower(col(col_name))))
```

Fig:3.4 Lowering & Trimming

Replace Empty Strings with Nulls

The empty strings in the dataset function as missing values because they are substituted with null which helps maintain data validity from the fig 3.5.

```

1 # Replace empty strings with nulls
2 from pyspark.sql.functions import when
3
4 for c in df.columns:
5     df = df.withColumn(c, when(col(c) == "", None).otherwise(col(c)))

```

Fig:3.5 Replacing Strings

Fill Missing Values

The algorithm fills numeric values missing with the input value of zero. The procedure for handling missing values within string columns involves substituting them with "unknown" entries from the fig 3.6. The completion of the data set and prevention of training-stage errors or transformations make this approach essential.

```

1 # Fill missing numeric values with 0 or median (simplified here as 0)
2 num_cols = [col_name for col_name, dtype in df.dtypes if dtype in ['int', 'double']]
3 df = df.fillna(0, subset=num_cols)
4 # Fill string columns with "unknown"
5 df = df.fillna("unknown", subset=string_columns)

```

Fig:3.6 Filling Missing Values

Step 4: Split Dataset by Property Purpose

To tailor the prediction models to specific property purposes, the dataset is divided into two subsets:

- **buyer_df:** The dataset exists for predicting the owner price per square foot. The feature set contains no attributes associated with property rentals because it omits rent_monthly_cost and Deposit Cost information.
- **rental_df:** The target variable in this dataset stands for predicting monthly rent amounts. Owner Price per Sqft and GOV Price per Sqft together with every feature related to buying were removed from this dataset from the fig 3.7.

```

1 # Step 3: Partition
2 buyer_df = df.filter(col("Property Purpose") == "buy").drop("rent_montly_cost", "Deposit Cost")
3 rental_df = df.filter(col("Property Purpose") == "rent").drop("Owner Price per Sqft", "GOV Price per Sqft")
4

```

Fig:3.7 Splitting Dataset

Step 5: Buyer Price Prediction Model

Feature Engineering

To enhance the predictive power of our model, we engineer new features from the fig 3.8:

- **BHK_Sqft**: Calculated by multiplying the number of BHK (bedrooms) with the square footage.
- **Months_To_Availability**: The number of months until the property becomes available, calculated using the Availability Date.
- **Has_Photos**: A binary feature indicating whether the property has photos or not.
- **Amenity_Score**: A score aggregating the presence of amenities like a gym, swimming pool, security, and convention hall.
- **Is_Near_IT_Hub**: A binary feature indicating whether the property is located near an IT hub.

```

1 from pyspark.sql.functions import col, months_between, current_date, when, lower
2
3 # 1. BHK x Sqft
4 ml_df = buyer_df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
5
6 # 2. Months to Availability
7 ml_df = ml_df.withColumn("Months_To_Availability", months_between(col("Availability Date"), current_date()))
8
9 # 3. Has Photos
10 ml_df = ml_df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
11
12 # 4. Amenity Score (gym, pool, security, convention hall)
13 ml_df = ml_df.withColumn("Amenity_Score",
14     when(lower(col("Gym")) == "yes", 1).otherwise(0) +
15     when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
16     when(lower(col("Security")) != "none", 1).otherwise(0) +
17     when(lower(col("Convention Hall")) == "yes", 1).otherwise(0))
18 )
19
20 # 5. Is Near IT Hub
21 ml_df = ml_df.withColumn("Is_Near_IT_Hub", when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
22

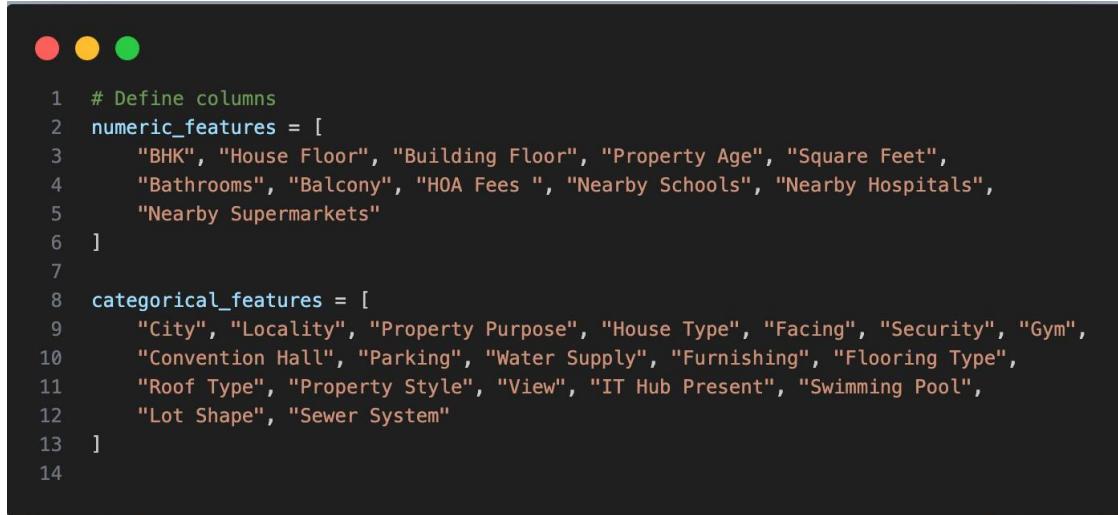
```

Fig:3.8 Feature Engineering

These new features are added to the dataset, improving its richness for prediction.

Feature Columns

- **Numeric Features:** Features that are continuous or numeric, such as BHK, Square Feet, Property Age, and Bathrooms.
- **Categorical Features:** Features that are categorical, such as City, Locality, Security, and Gym.



```
● ● ●  
1 # Define columns  
2 numeric_features = [  
3     "BHK", "House Floor", "Building Floor", "Property Age", "Square Feet",  
4     "Bathrooms", "Balcony", "HOA Fees ", "Nearby Schools", "Nearby Hospitals",  
5     "Nearby Supermarkets"  
6 ]  
7  
8 categorical_features = [  
9     "City", "Locality", "Property Purpose", "House Type", "Facing", "Security", "Gym",  
10    "Convention Hall", "Parking", "Water Supply", "Furnishing", "Flooring Type",  
11    "Roof Type", "Property Style", "View", "IT Hub Present", "Swimming Pool",  
12    "Lot Shape", "Sewer System"  
13 ]  
14
```

Fig:3.9 Feature Columns

We separate the numeric and categorical features to handle them differently during model training as shown in the fig 3.8.

Step 6: ML Pipeline for Encoding & Vector Assembly

We built a pipeline from fig 3.10.

- **String Indexing:** Converting categorical columns into numeric indices.
- **One-Hot Encoding:** Encoding categorical indices into binary vectors.
- **Vector Assembly:** Combining all features (numeric and one-hot encoded) into a single vector column (features) for model input.

```

1 # Index and OneHotEncode categorical features
2 indexers = [StringIndexer(inputCol=col, outputCol=col + "_Index", handleInvalid='keep') for col in categorical_features]
3 encoders = [OneHotEncoder(inputCol=col + "_Index", outputCol=col + "_OHE") for col in categorical_features]
4
5 # Final input features (numeric + encoded)
6 encoded_features = [col + "_OHE" for col in categorical_features]
7 assembler_inputs = numeric_features + encoded_features
8
9 assembler = VectorAssembler(
10     inputCols=assembler_inputs,
11     outputCol="features"
12 )
13
14 # Target column
15 target_column = "Owner Price per Sqft"
16
17 pipeline = Pipeline(stages=indexers + encoders + [assembler])
18 pipeline_model = pipeline.fit(ml_df)
19 processed_df = pipeline_model.transform(ml_df)

```

Fig:3.10 ML Pipeline

This pipeline prepares the dataset for machine learning by transforming it into the required format for Spark MLlib.

Step 7:Train-Test Split and Model Training

We split the dataset into training and testing sets (80% training, 20% testing). Then, we train a **Linear Regression** model to predict the Owner Price per Sqft as shown in the fig 3.11.

```

1 train_df, test_df = processed_df.randomSplit([0.8, 0.2], seed=42)
2 from pyspark.ml.regression import LinearRegression
3
4 lr = LinearRegression(
5     featuresCol="features",
6     labelCol=target_column,
7     predictionCol="prediction"
8 )
9
10 lr_model = lr.fit(train_df)
11 results = lr_model.evaluate(test_df)
12
13 print("RMSE:", results.rootMeanSquaredError)
14 print("R2:", results.r2)

```

Fig:3.11 Testing and Training

Model evaluation metrics include :

- **RMSE (Root Mean Squared Error):** ≈ 3188
- **R²:** ≈ 0.82 (Good predictive performance)

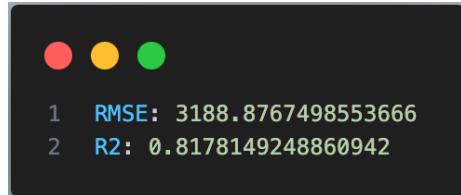


Fig:3.12 Model Evaluation Metrics

These results indicate that the model can predict the price with reasonable accuracy.

Step 8: Price Prediction for New Input

We can use the trained model to predict property prices for new data from the fig 3.13:

1. **User Input:** We create a user input row based on the property details.
2. **Feature Engineering:** Apply the same transformations to the user input as done for the training data.
3. **Prediction:** The model predicts the price per square foot, which is then used to compute the total property price.

```
1 from pyspark.sql.functions import round as spark_round
2
3 # 2. Add Total_Predicted_Price column
4 prediction = prediction.withColumn(
5     "Total_Predicted_Price",
6     spark_round(col("prediction") * col("Square Feet"))
7 )
8
9 # 3. Get values from prediction DataFrame
10 row = prediction.select("prediction", "Square Feet", "Total_Predicted_Price").first()
11
12 # 4. Extract values
13 predicted_price_per_sqft = round(row["prediction"])
14 total_price = row["Total_Predicted_Price"]
15
16 # 🎉 Display nicely
17 print(f"➡ Predicted Price per Sqft: ₹{predicted_price_per_sqft:,}")
18 print(f"🏡 Total Estimated Property Price: ₹{total_price:,}")
```

Fig:3.13 Price Prediction

Example output:

- **Predicted Price per Sqft:** ₹ [predicted value]
- **Total Estimated Property Price:** ₹ [predicted value]

This approach allows us to easily estimate property prices for new properties by leveraging the trained model from the table 3.4.

Summary Table			
Task	Target	RMSE	R ²
Price Prediction (Buy)	Owner Price per Sqft	3188.87	0.8178
Rent Prediction	Monthly Rent	11004.18	0.5816

Table:3.4 Model Evaluation Metrics

3.5 Customer Page

The Streamlit-based web application functions as a user interface which provides real estate listing search and Google Gemini AI recommendations following user-defined filter parameters. The web application obtains data from AWS S3 storage to achieve scalability in practical usage from the fig 3.14.

Libraries Used:

- streamlit: For building the web UI.
- boto3: AWS SDK to access S3.
- pandas: For data handling and filtering.
- openai, genai: To interface with OpenAI or Gemini (only Gemini is used here).
- datetime, os: Utility modules.

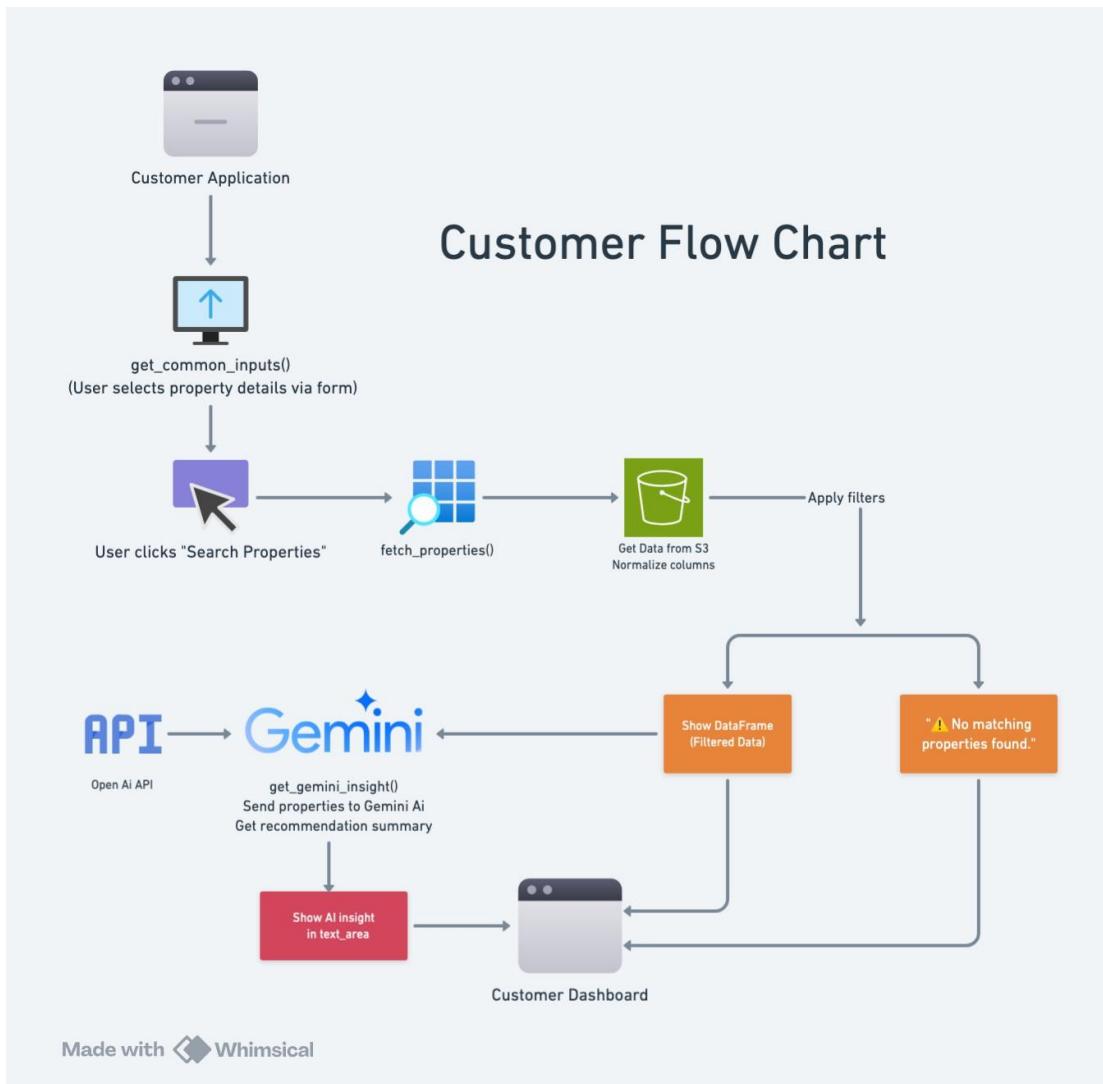


Fig:3.14 Customer flow chart

3.5.1 Gemini Setup:

```
client = genai.Client(api_key=...) # Sets up Gemini AI API
```

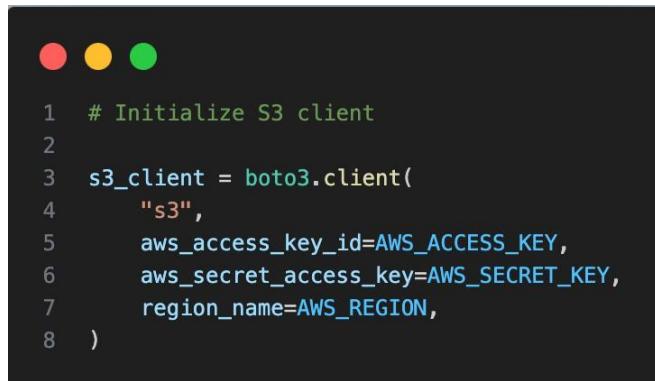
```

● ● ●
1 client = genai.Client(api_key='AIzaSyBY4Rj0QyMtK641PZBrwFTaqwRInHmqzEk')
  
```

Fig:3.15 Gemini Api Key

3.5.2 AWS S3 Setup: from the fig 3.16

```
s3_client = boto3.client("s3", ...)
```



```
● ● ●
1 # Initialize S3 client
2
3 s3_client = boto3.client(
4     "s3",
5     aws_access_key_id=AWS_ACCESS_KEY,
6     aws_secret_access_key=AWS_SECRET_KEY,
7     region_name=AWS_REGION,
8 )
```

Fig:3.16 AWS S3 setup

Used to fetch CSV datasets from an S3 bucket.

3.5.3 get_common_inputs() Function

This creates the **user input form** in the sidebar/UI using Streamlite widgets like select box, slider, number_input, etc.

It collects property features like:

- Purpose: Buy/Rent
- Location: City & Locality
- House Features: BHK, Square feet, Floor, Age
- Amenities: Gym, Pool, Security
- Proximity: Schools, Hospitals, Supermarkets
- Availability, Style, Furnishing, etc.

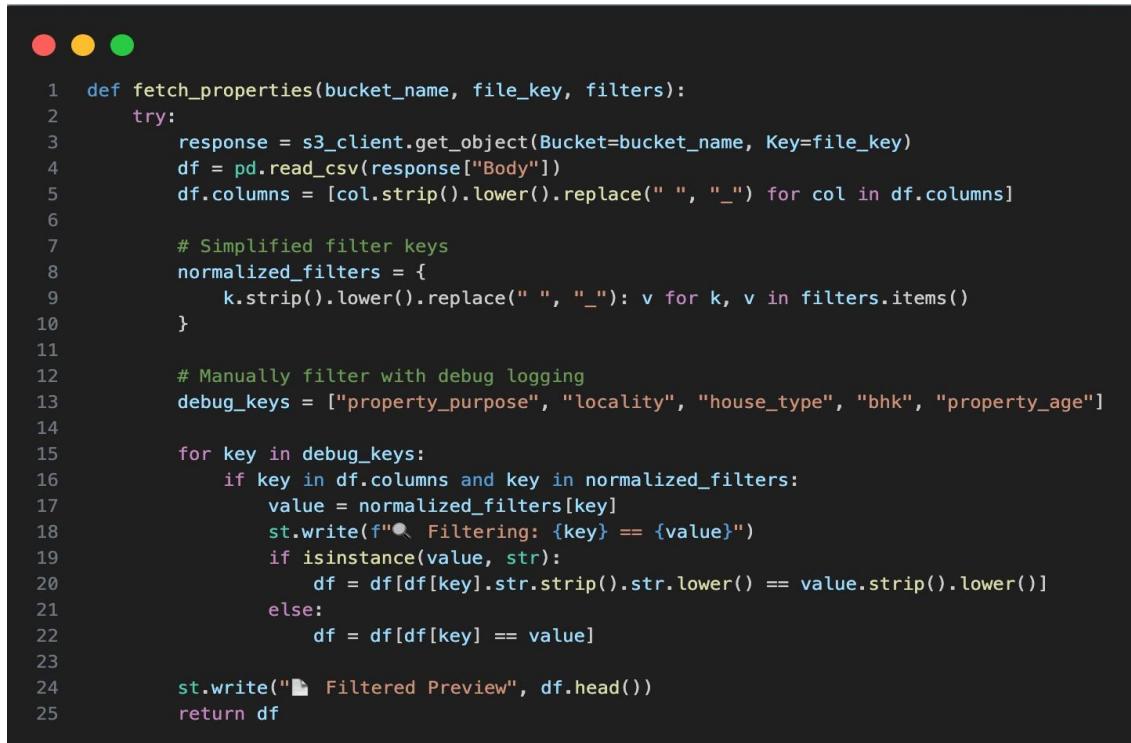
All inputs are returned in a structured dictionary. This serves as a **filter** for querying the dataset.

3.5.4 fetch_properties() Function

This fetches the dataset from **AWS S3**, processes it with Pandas, and **filters** based on the user input from the fig 3.17.

Key Steps:

- Downloads CSV (real_state_data.csv) from S3 bucket mage dataset.
- Normalizes column names for matching.
- Applies filters (like property_purpose, locality, BHK) and returns a filtered DataFrame.



```
1 def fetch_properties(bucket_name, file_key, filters):
2     try:
3         response = s3_client.get_object(Bucket=bucket_name, Key=file_key)
4         df = pd.read_csv(response["Body"])
5         df.columns = [col.strip().lower().replace(" ", "_") for col in df.columns]
6
7         # Simplified filter keys
8         normalized_filters = {
9             k.strip().lower().replace(" ", "_"): v for k, v in filters.items()
10            }
11
12         # Manually filter with debug logging
13         debug_keys = ["property_purpose", "locality", "house_type", "bhk", "property_age"]
14
15         for key in debug_keys:
16             if key in df.columns and key in normalized_filters:
17                 value = normalized_filters[key]
18                 st.write(f"● Filtering: {key} == {value}")
19                 if isinstance(value, str):
20                     df = df[df[key].str.strip().str.lower() == value.strip().lower()]
21                 else:
22                     df = df[df[key] == value]
23
24         st.write("Filtered Preview", df.head())
25     return df
```

Fig:3.17 Fetching Properties

It also displays the filtered preview using st.write().

3.5.5 get_gemini_insight() Function

This function uses the **Google Gemini API** to:

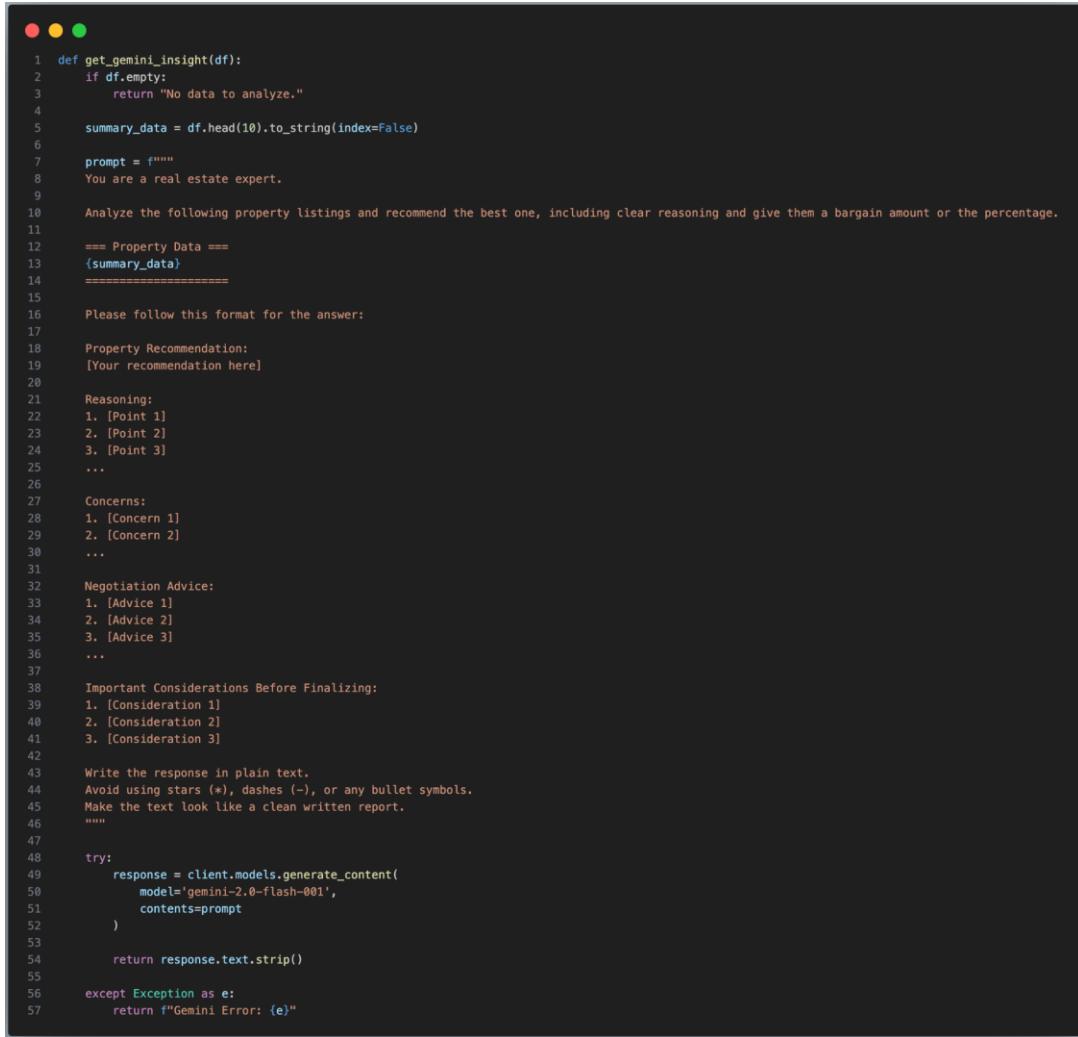
- Analyze the top 10 filtered property listings.
- Generate a **recommendation** report.

The prompt sent to Gemini includes from the fig 3.18:

- Summary of listings.
- Instructions to return:

- Best property recommendation.

- Reasoning
- Concerns
- Negotiation advice
- Final considerations



```
1 def get_gemini_insight(df):
2     if df.empty:
3         return "No data to analyze."
4
5     summary_data = df.head(10).to_string(index=False)
6
7     prompt = f"""
8     You are a real estate expert.
9
10    Analyze the following property listings and recommend the best one, including clear reasoning and give them a bargain amount or the percentage.
11
12    === Property Data ===
13    {summary_data}
14    =====
15
16    Please follow this format for the answer:
17
18    Property Recommendation:
19    [Your recommendation here]
20
21    Reasoning:
22    1. [Point 1]
23    2. [Point 2]
24    3. [Point 3]
25    ...
26
27    Concerns:
28    1. [Concern 1]
29    2. [Concern 2]
30    ...
31
32    Negotiation Advice:
33    1. [Advice 1]
34    2. [Advice 2]
35    3. [Advice 3]
36    ...
37
38    Important Considerations Before Finalizing:
39    1. [Consideration 1]
40    2. [Consideration 2]
41    3. [Consideration 3]
42
43    Write the response in plain text.
44    Avoid using stars (*), dashes (-), or any bullet symbols.
45    Make the text look like a clean written report.
46    """
47
48    try:
49        response = client.models.generate_content(
50            model='gemini-2.0-flash-001',
51            contents=prompt
52        )
53
54        return response.text.strip()
55
56    except Exception as e:
57        return f"Gemini Error: {e}"
```

Fig:3.18 Plain Text Report

Gemini returns this as a plain-text report, which is displayed in a text area.

3.5.6 main() Function – The App Entry Point

This runs the actual app logic from the fig 3.19.

Steps:

1. Title: “Real Estate Property Insights for Customers”
2. Collects input using get_common_inputs()
3. On button click:
 - Fetches property listings with filters from S3
 - Shows results in a table
 - Sends results to Gemini for recommendation
 - Displays AI-generated recommendation



```
1 # 🚀 Main App
2 def main():
3     st.title("🏡 Real Estate Property Insights for Customers")
4     bucket_name = "magedataset"
5     file_key = "real_state_data.csv"
6
7     filters = get_common_inputs(is_rent=True)
8
9     if st.button("Search Properties"):
10         result_df = fetch_properties(bucket_name, file_key, filters)
11         if not result_df.empty:
12             st.success(f"✅ Found {len(result_df)} matching properties.")
13             st.dataframe(result_df)
14             insight = get_gemini_insight(result_df)
15             st.text_area("🧠 OpenAI Recommendation", insight, height=200)
16         else:
17             st.warning("⚠️ No matching properties found.")
18
19     if __name__ == "__main__":
20         main()
21
```

Fig:3.19 Main Function Setup

3.6 Owner

The AI-powered tool named PropIntel provides genuine real estate guidance across properties in Chennai to home buyers along with property owners. The real-time AI insights along with advanced machine learning models which PropIntel provides help users find properties for buying or renting from the fig 3.20.

Accurate price or rent predictions based on your property details

Smart, personalized market advice powered by OpenAI's GPT

An intuitive and interactive interface to guide your real estate journey

3.6.1 Libraries and Initialization

- Streamlit (st): For building the web interface.
- pandas: Standard data manipulation (though not heavily used here).
- openai & google.genai: For generating AI-based insights.
- Spark (pySpark): Handles large-scale data processing and prediction.
- Pre-trained ML models: Used for price/rent prediction.

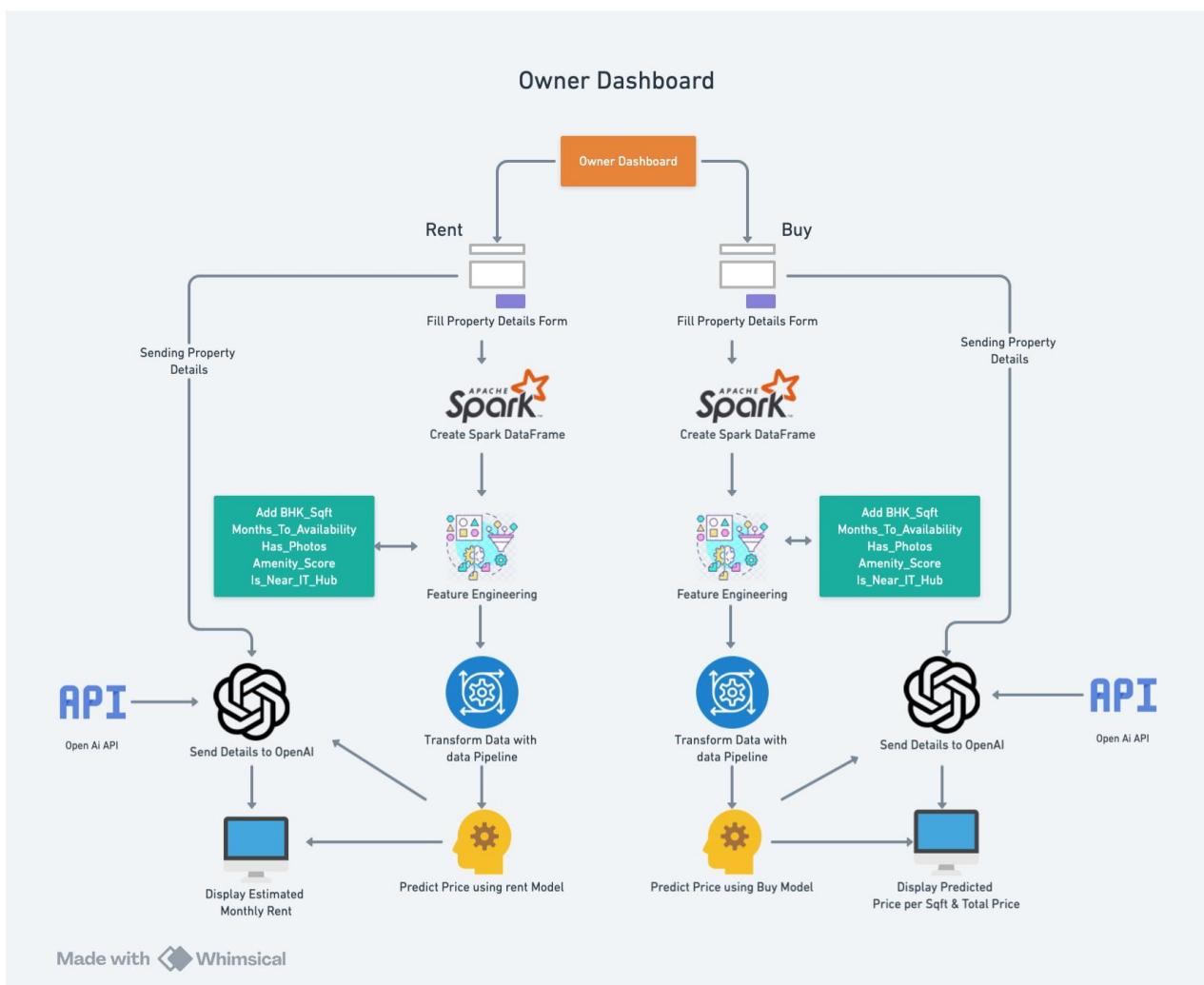
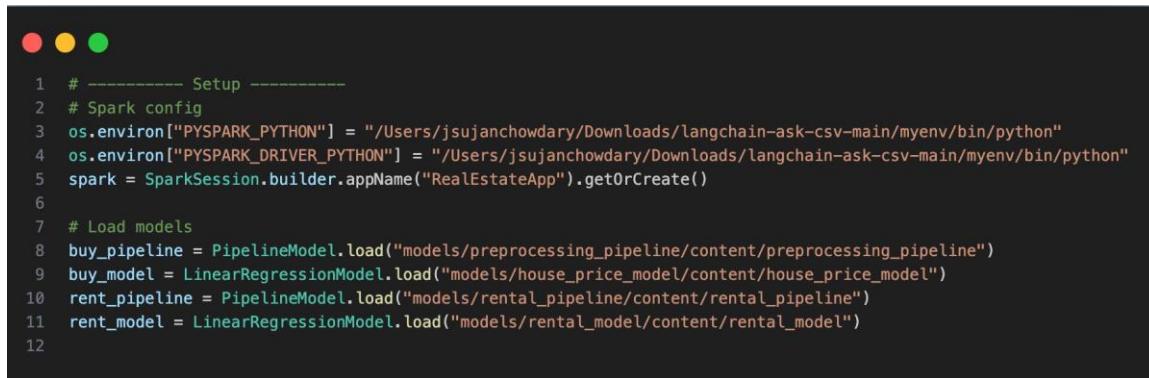


Fig:3.20 Owner Dashboard

3.6.2 Model Loading

Spark ML models are loaded:

- buy_pipeline, buy_model for buy price prediction.
- rent_pipeline, rent_model for rental price prediction.



```
1 # ----- Setup -----
2 # Spark config
3 os.environ["PYSPARK_PYTHON"] = "/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/myenv/bin/python"
4 os.environ["PYSPARK_DRIVER_PYTHON"] = "/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/myenv/bin/python"
5 spark = SparkSession.builder.appName("RealEstateApp").getOrCreate()
6
7 # Load models
8 buy_pipeline = PipelineModel.load("models/preprocessing_pipeline/content/preprocessing_pipeline")
9 buy_model = LinearRegressionModel.load("models/house_price_model/content/house_price_model")
10 rent_pipeline = PipelineModel.load("models/rental_pipeline/content/rental_pipeline")
11 rent_model = LinearRegressionModel.load("models/rental_model/content/rental_model")
12
```

Fig:3.21 Model Loading

These pipelines likely include feature preprocessing steps like encoding, scaling, etc.

3.6.3 User Interface (UI)

- Built using Streamlit.
- User chooses between:
Buy House Price

Rent Monthly Cost

Then, they fill in a form about their property — location, type, amenities, age, etc.

3.6.4 Form Input Collection

Function: get_common_inputs(is_rent=False)

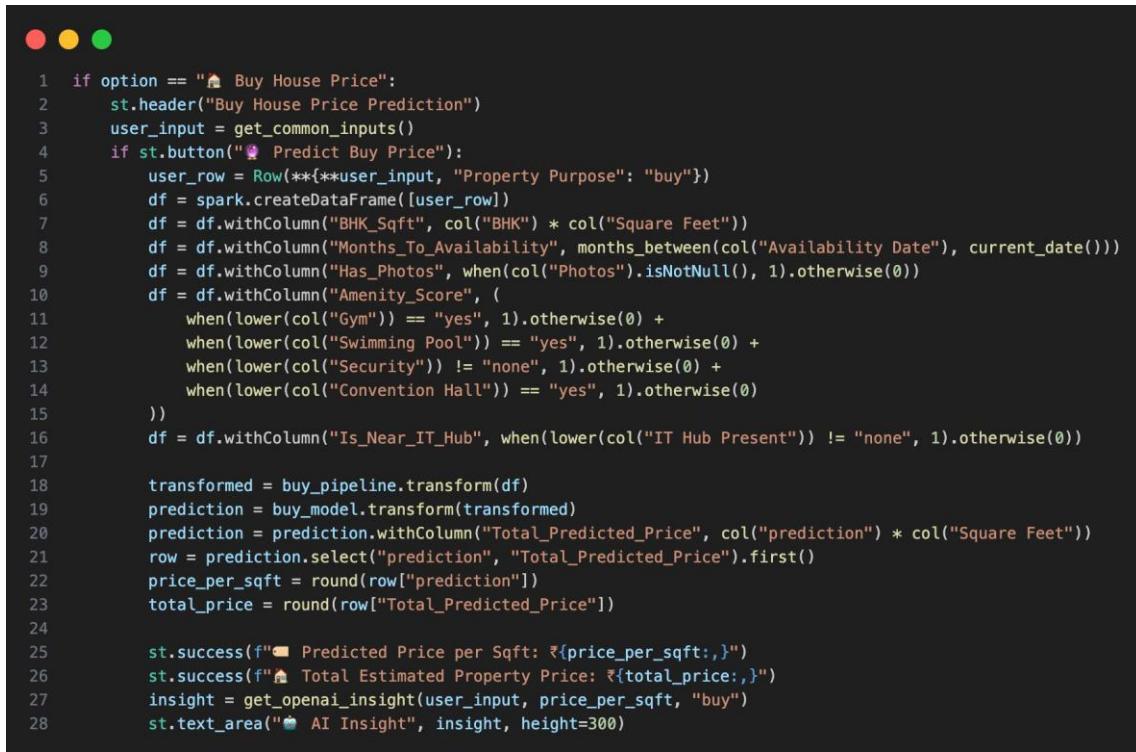
- Collects all user inputs.
- For rent predictions, extra fields like Deposit and separate HOA fees are added.

Returns a Python dict with all input values.

3.6.5 Predictions via Spark ML Models

Depending on the user's choice:

- The input is converted into a Spark DataFrame.
- Feature engineering is applied:
BHK_Sqft, Months_To_Availability, Has_Photos, etc from the fig 3.22.
- Transformed via pipeline.
- Prediction made using the corresponding model (buy_model or rent_model).
- Result is displayed on the UI.



```
1 if option == "🏡 Buy House Price":
2     st.header("Buy House Price Prediction")
3     user_input = get_common_inputs()
4     if st.button("🔮 Predict Buy Price"):
5         user_row = Row(**{*user_input, "Property Purpose": "buy"})
6         df = spark.createDataFrame([user_row])
7         df = df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
8         df = df.withColumn("Months_To_Availability", months_between(col("Availability Date"), current_date()))
9         df = df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
10        df = df.withColumn("Amenity_Score", (
11            when(lower(col("Gym")) == "yes", 1).otherwise(0) +
12            when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
13            when(lower(col("Security")) != "none", 1).otherwise(0) +
14            when(lower(col("Convention Hall")) == "yes", 1).otherwise(0)
15        ))
16        df = df.withColumn("Is_Near_IT_Hub", when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
17
18        transformed = buy_pipeline.transform(df)
19        prediction = buy_model.transform(transformed)
20        prediction = prediction.withColumn("Total_Predicted_Price", col("prediction") * col("Square Feet"))
21        row = prediction.select("prediction", "Total_Predicted_Price").first()
22        price_per_sqft = round(row["prediction"])
23        total_price = round(row["Total_Predicted_Price"])
24
25        st.success(f"➡️ Predicted Price per Sqft: ₹{price_per_sqft:,}")
26        st.success(f"🏡 Total Estimated Property Price: ₹{total_price:,}")
27        insight = get_openai_insight(user_input, price_per_sqft, "buy")
28        st.text_area("🌐 AI Insight", insight, height=300)
```

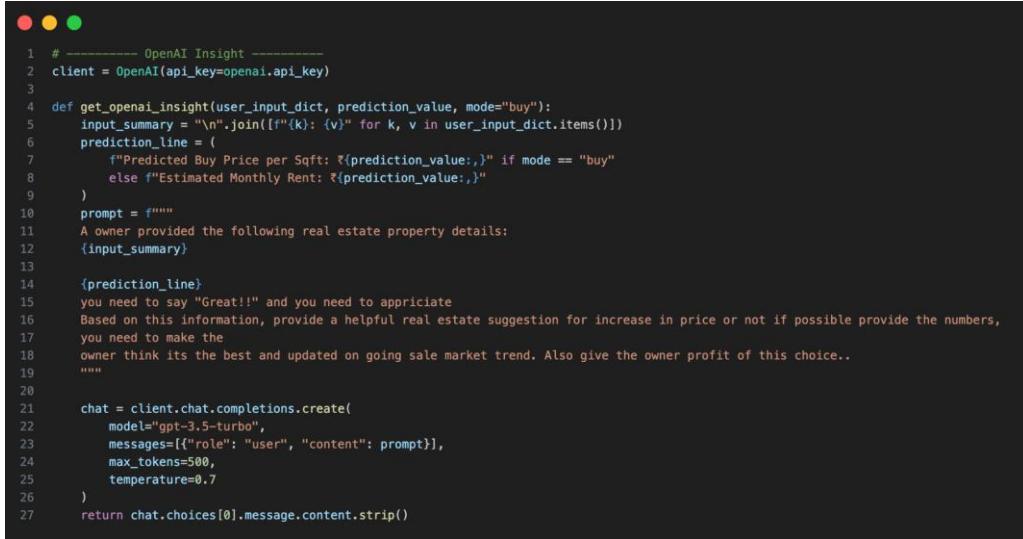
Fig:3.22 Prediction via Spark ML

3.6.6 OpenAI Chat Insight

Function: `get_openai_insight(...)`

- Builds a natural language summary of the input and predicted value.
- Sends this to GPT-3.5 using openai.ChatCompletion.
- Provides a personalized message to the user from the fig 3.23.

- Encouragement ("Great!!")
- Trends and insights about the market
- Possibly estimated profits or future value suggestion



```

● ● ●
1 # ----- OpenAI Insight -----
2 client = OpenAI(api_key=openai.api_key)
3
4 def get_openai_insight(user_input_dict, prediction_value, mode="buy"):
5     input_summary = "\n".join([f"{k}: {v}" for k, v in user_input_dict.items()])
6     prediction_line = (
7         f"Predicted Buy Price per Sqft: ${prediction_value}, if mode == \"buy\""
8         else f"Estimated Monthly Rent: ${prediction_value},"
9     )
10    prompt = f"""
11    A owner provided the following real estate property details:
12    {input_summary}
13
14    {prediction_line}
15    you need to say "Great!!" and you need to appriciate
16    Based on this information, provide a helpful real estate suggestion for increase in price or not if possible provide the numbers,
17    you need to make the
18    owner think its the best and updated on going sale market trend. Also give the owner profit of this choice..
19    """
20
21    chat = client.chat.completions.create(
22        model="gpt-3.5-turbo",
23        messages=[{"role": "user", "content": prompt}],
24        max_tokens=500,
25        temperature=0.7
26    )
27    return chat.choices[0].message.content.strip()

```

Fig:3.23 Open AI Chat Insight

3.6.7 Sample Output

If the user fills in details and clicks "Predict" from the fig 3.24.

- The app shows:
 - Price per sqft and total property value (or rent)
 - AI-generated insight in a text area

```
1 # ----- Rent Prediction -----
2 elif option == "₹ Rent Monthly Cost":
3     st.header("Rent Monthly Cost Prediction")
4     user_input = get_common_inputs(is_rent=True)
5     if st.button("Predict Rent Price"):
6         user_row = Row(**{**user_input, "Property Purpose": "rent"})
7         df = spark.createDataFrame([user_row])
8         df = df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
9         df = df.withColumn("Months_To_Availability", months_between(col("Availability Date"), current_date()))
10        df = df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
11        df = df.withColumn("Amenity_Score", (
12            when(lower(col("Gym")) == "yes", 1).otherwise(0) +
13            when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
14            when(lower(col("Security")) != "none", 1).otherwise(0) +
15            when(lower(col("Convention Hall")) == "yes", 1).otherwise(0)
16        ))
17        df = df.withColumn("Is_Near_IT_Hub", when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
18
19        transformed = rent_pipeline.transform(df)
20        prediction = rent_model.transform(transformed)
21        rent = round(prediction.select("prediction").first()["prediction"])
22
23        st.success(f"₹ Estimated Monthly Rent: ₹{rent:,}")
24        insight = get_openai_insight(user_input, rent, "rent")
25        st.text_area("AI Insight", insight, height=300)
```

Fig:3.24 Sample Output

CHAPTER-4

IMPLEMENTATION OF REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM

4.1 System Design

The system is designed as a modular real estate prediction platform, integrating user interaction through a Streamlit web app, backend computation with Apache Spark, and AI insights via OpenAI's GPT model. The high-level design is divided into the following modules:

- Frontend UI (Streamlit): Captures user inputs and displays predictions/insights.
- Feature Engineering and Preprocessing (PySpark): Transforms raw input into features suitable for ML models.
- Machine Learning Models (Spark MLlib): Trained regression models for price and rent prediction.
- AI Insights Engine (OpenAI GPT): Generates real estate advice based on prediction results and property features.

4.2 Coding and Tools

- **Languages:** Python
- **Libraries and Frameworks:**
 - Streamlit: For building the web interface
 - PySpark: For data processing and model serving
 - pandas: Lightweight data handling
 - openai: API client for GPT insights
 - google.generativeai: (Imported but not used — may be for future GenAI extensions)
- **Development Environment:**
 - Local Python virtual environment
 - Pre-trained Spark ML models stored and loaded from local models/ directory

4.3 Process Flow

The process is bifurcated into two main flows based on user selection: Buy House Price or Rent Monthly Cost. Here's the generalized flow from table 4.1.

Common Flow:

1. **User Input Collection:** Users input property details via the Streamlit interface.
2. **Spark DataFrame Conversion:** Inputs are wrapped in a Spark Row and converted to a Spark DataFrame.
3. **Feature Engineering:**
 - Create derived columns: BHK_Sqft, Months_To_Availability, Has_Photos, Amenity_Score, etc.
 - Encode categorical features.
4. **Model Inference:**
 - For buying, use buy_pipeline and buy_model.
 - For renting, use rent_pipeline and rent_model.
5. **Result Interpretation:**
 - For buying: Calculate price per square foot and total estimated price.
 - For renting: Output monthly rent cost.
6. **OpenAI Integration:**
 - Construct a prompt using user inputs and prediction results.
 - Call OpenAI's Chat API to return a natural-language recommendation and market insight.
7. **UI Output:**
 - Display prediction results and AI insights using Streamlit's st.success() and st.text_area().

Buy Flow:

- Uses models/house_price_model for predicting price per square foot.
- Total price is calculated using predicted price per square foot \times input square footage.

Rent Flow:

- Uses models/rental_model for predicting monthly rent directly.

4.4 Challenges and Obstacles

Challenge	Description
Model Integration with PySpark	Converting user input to a format compatible with Spark ML models, especially when using Row and applying transformations dynamically.
Handling Categorical Inputs	Ensuring proper encoding and mapping of string-based categorical data to what the model expects (e.g., locality, amenities, style).
Date Calculations in Spark	Calculating Months_To_Availability from the Availability Date using Spark's date functions while maintaining correct format.
OpenAI Rate Limit & Latency	Balancing responsiveness and prompt quality while staying within API limits and avoiding long waits for responses.
Streamlit Performance	Ensuring the UI remains smooth despite backend Spark computations, especially in local dev environments.
Model Version Management	Managing and loading correct versions of PipelineModel and LinearRegressionModel from disk consistently.

Table:4.1 Challenges and Obstacles

4.5 Cloud Deployment and Integration (AWS + GCP)

To enable scalable, efficient, and cloud-native deployment of the real estate prediction platform, we implemented a hybrid multi cloud architecture using both Amazon Web Services (AWS) and Google Cloud Platform (GCP). This allowed us to leverage the strengths of each platform for different parts of the ML lifecycle—data storage, model training, visualization, and deployment.

4.5.1 AWS Components

4.5.1.1 Amazon S3 – Centralized Data Storage

- Used as the single source of truth for all raw, processed, and model data from the fig 4.1
- Stores:
 - User input records and metadata
 - Trained PySpark PipelineModel and LinearRegressionModel
 - Feature-engineered datasets
 - Prediction outputs

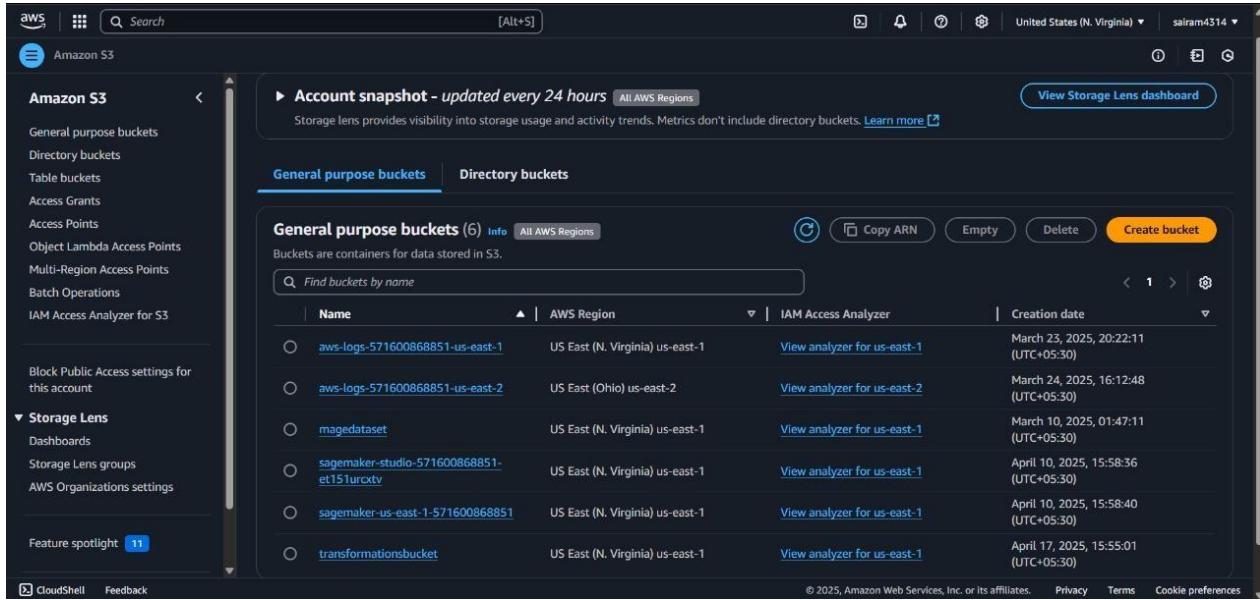


Fig:4.1 AWS S3 Setup

4.5.1.2 Amazon SageMaker Studio – Model Visualization & Experimentation

- SageMaker Studio was used to:
 - Access and load model outputs from S3.
 - Perform exploratory data analysis (EDA) and validation on predictions.
 - Create quick transformations for real estate rent and buy datasets.
- SageMaker notebooks enabled us to experiment with alternate ML pipelines in a managed environment from the fig 4.2.

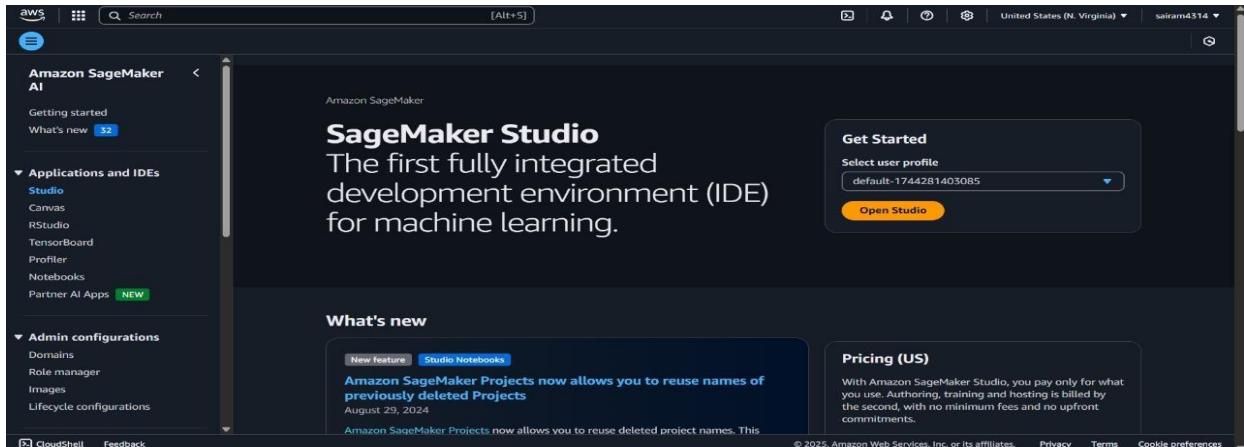


Fig:4.2 SageMaker Studio Setup

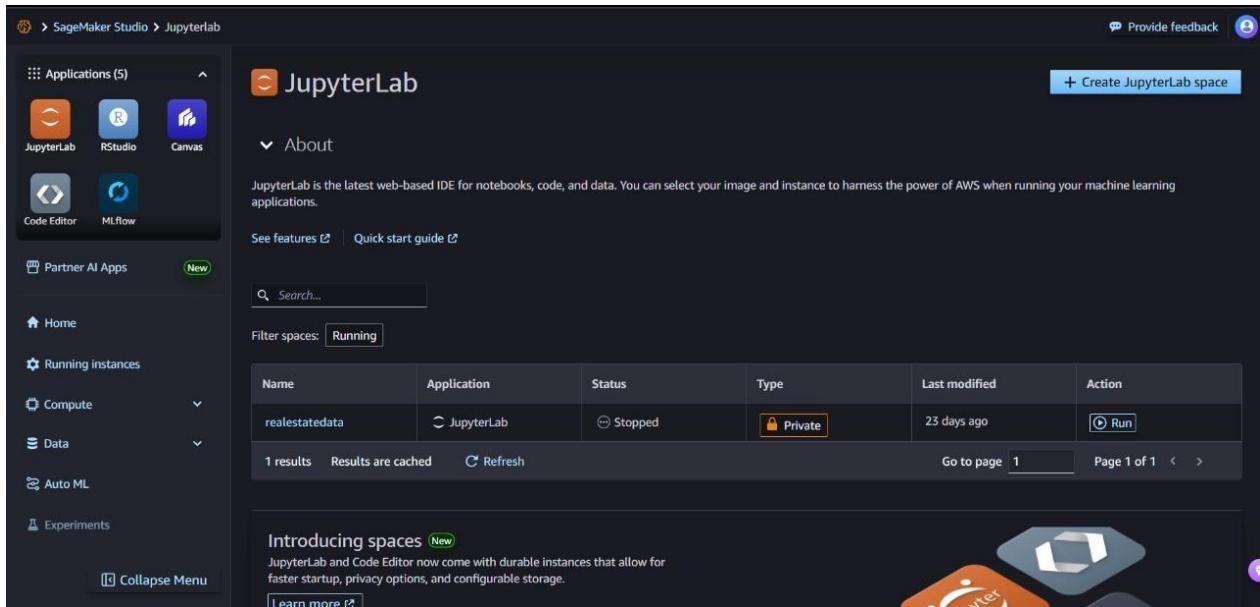


Fig:4.3 JupyterLab in Sage maker Studio

The screenshot shows the SageMaker Jupyter Lab interface. The top menu includes File, Edit, View, Run, Kernel, Git, Tabs, Settings, and Help. The tabs show two notebooks: 'Untitled1.ipynb' and 'Untitled2.ipynb'. The code editor displays the following Python code:

```

[1]: pip install pyspark
Requirement already satisfied: pyspark in /opt/conda/lib/python3.11/site-packages (3.5.5)
Requirement already satisfied: py4j==0.10.9.7 in /opt/conda/lib/python3.11/site-packages (from pyspark) (0.10.9.7)
Note: you may need to restart the kernel to use updated packages.

[2]: pip install boto3
Requirement already satisfied: boto3 in /opt/conda/lib/python3.11/site-packages (1.36.23)
Requirement already satisfied: botocore<1.37.0,>=1.36.23 in /opt/conda/lib/python3.11/site-packages (from boto3) (1.36.23)
Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in /opt/conda/lib/python3.11/site-packages (from boto3) (1.0.1)
Requirement already satisfied: s3transfer<0.12.0,>=0.11.0 in /opt/conda/lib/python3.11/site-packages (from boto3) (0.11.3)
Requirement already satisfied: python-dateutil<3.0.0,>=2.1 in /opt/conda/lib/python3.11/site-packages (from botocore<1.37.0,>=1.36.23>;boto3) (2.9.0)
Requirement already satisfied: urllib3<2.2.0,>=1.25.4 in /opt/conda/lib/python3.11/site-packages (from botocore<1.37.0,>=1.36.23>;boto3) (2.3.0)
Requirement already satisfied: six<1.5 in /opt/conda/lib/python3.11/site-packages (from python-dateutil<3.0.0,>=2.1>;botocore<1.37.0,>=1.36.23>;boto3) (1.17.0)
Note: you may need to restart the kernel to use updated packages.

[3]: import os
import boto3

[4]: os.environ["AWS_ACCESS_KEY_ID"] = "AKIAVKFQRCXZX0XDHA66"
os.environ["AWS_SECRET_ACCESS_KEY"] = "pAHf8lHi/wtKg/DNUCJ3Ogc5juPhHirOK1tVzLue"

[5]: s3 = boto3.client('s3')
bucket = "magedataset"
key = "real_state_data.csv"
local_file = "/tmp/real_state_data.csv"

# Download the file
s3.download_file(bucket, key, local_file)
print("Download complete.")
Download complete.

```

At the bottom, there are status indicators for 'Simple', 'Instance MEM 28%', 'Amazon Q', and 'Cookie Preferences'.

Fig:4.4 SageMaker Jupyter Lab HomePage

4.5.1.3 Amazon QuickSight – Business Intelligence Dashboard

- Built interactive dashboards to visualize:
 - Predicted vs actual prices for different localities
 - Rent and buy price trends over time

- Amenity score impact on pricing
- **Data source:** Amazon S3 (auto-refreshed).
 - Helped stakeholders derive actionable insights from prediction trends from the fig 4.5.



Fig:4.5 Quicksight Buyer Dashboard

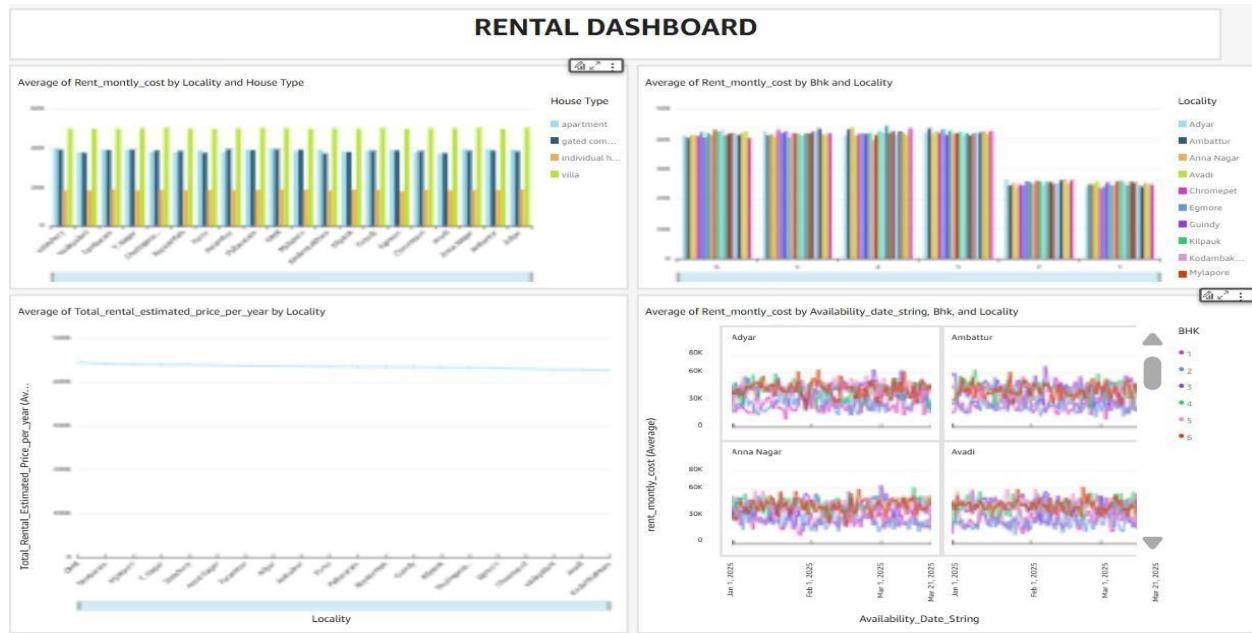


Fig:4.6 QuickSight Rental Dashboard

4.5.2 GCP Components

4.5.2.1 Vertex AI – Model Development and Deployment

- Used for:
 - Training and tuning machine learning models with integrated notebook support.
 - Running alternate versions of models for experimentation.
 - Creating endpoints for model predictions.
- Integrated directly with S3 using a multi cloud approach, allowing seamless access to AWS-stored data inside GCP.

4.5.2.2 Vertex Pipelines

- If needed, could be used to orchestrate training pipelines from S3 → Preprocessing → Model Training → Vertex Endpoint Deployment from the fig 4.7 and 4.8.

The screenshot shows the Google Cloud Platform interface for Vertex AI Workbench. The left sidebar has sections for Tools (Dashboard, Model Garden, Pipelines), Notebooks (Colab Enterprise), and Vertex AI Studio (Overview, Create prompt, Media Studio, Stream realtime, Prompt gallery, Prompt management). The main area is titled 'Workbench' with tabs for 'Create New' and 'Refresh'. Below is the 'Instances' tab, which is selected. A sub-menu shows 'View: Instances' (selected), User-managed Notebooks, and Managed Notebooks. A message box says 'JupyterLab 4 is now available in Vertex AI Workbench.' A 'Dismiss' button is in the top right of the message box. The main content area displays a table of instances:

Instance name	Zone	Auto upgrade	Version	Machine Type	GPUs	Owner	Cr.
instance_20250420-120545	us-central1-a	—	M129	Efficient Instance: 4 vCPUs, 16 GB RAM	None	23193527377-compute@developer.gserviceaccount.com	Ap 20 12

Fig:4.7 Google Vertex AI

```

[1]: import os
import boto3

[2]: s3 = boto3.client('s3')
bucket = "mageddataset"
key = "real_state_data.csv"
local_file = "/tmp/real_state_data.csv"

# Download the file
s3.download_file(bucket, key, local_file)
print("Download complete.")

Download complete.

[3]: from pyspark.sql import SparkSession
spark = SparkSession.builder.appName("RealStateData").getOrCreate()
df = spark.read.csv(local_file, header=True, inferSchema=True)

df.printSchema()
df.show()

Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
25/04/20 07:19:20 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
root
|-- Property_ID: string (nullable = true)

```

Fig:4.8 Vertex AI HomePage

4.5.3 Multi Cloud Architecture: Why & How

Component	Cloud Provider	Reason
Data Storage (S3)	AWS	Durable, highly available, and flexible storage accessible from both clouds.
Model Experimentation	GCP Vertex AI	Better integration with custom ML pipelines and auto-ML options.
BI Dashboard	AWS QuickSight	Native AWS tool for fast business visualization over S3 datasets.

Table:4.2 Detailing of Components

4.5.4 Multi Cloud Strategy Highlights:

- Cost efficiency: Only necessary services are used from each platform.
- Flexibility: Experiment with models on GCP while storing all assets in S3.
- Interoperability: Vertex AI accessed S3 via service accounts with proper IAM permissions or pre-signed URLs.

4.5.5 Security & Access Control

- **IAM Roles and Policies:** Applied least privilege policies for both SageMaker and Vertex AI to access only required S3 buckets from the fig 4.9.
- **API Key Management:** OpenAI and internal secrets are stored securely using AWS Secrets Manager or Vertex AI's secret integration.

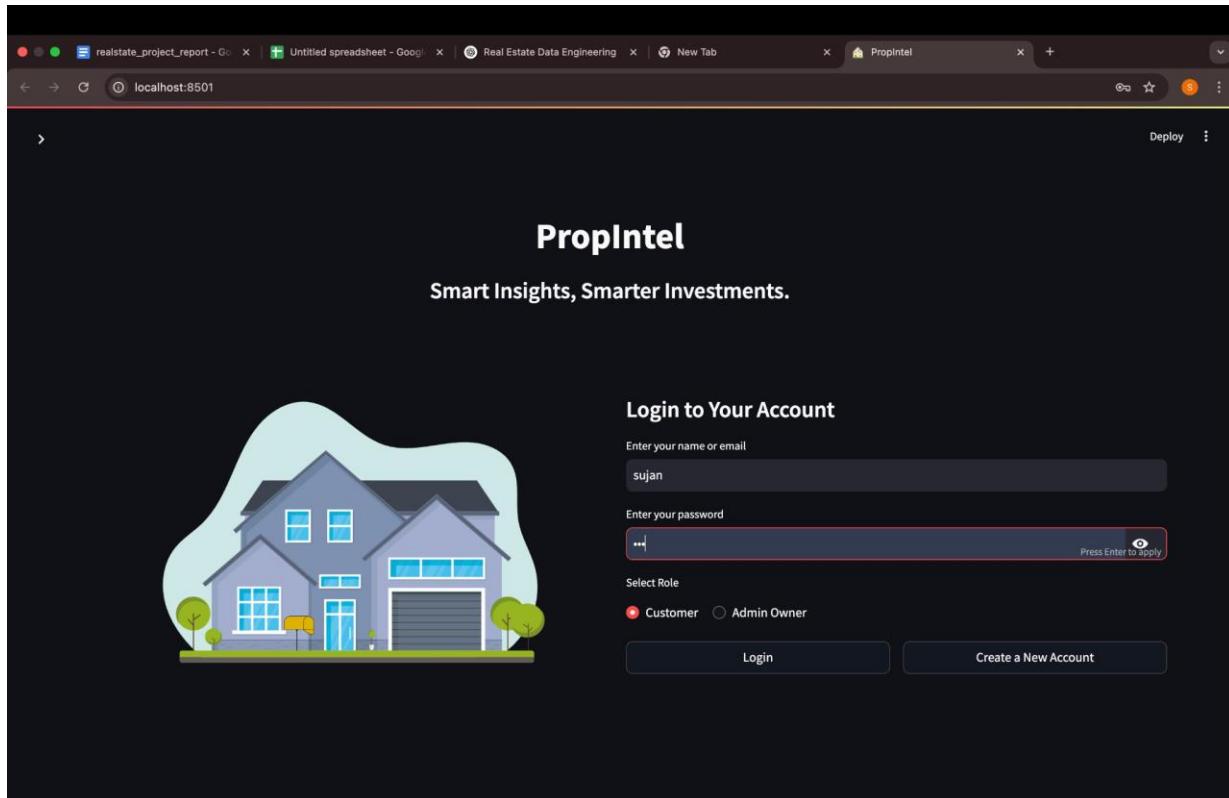


Fig:4.9 Login page

CHAPTER 5

RESULTS AND DISCUSSIONS REAL ESTATE DATA ENGINEERING USING MULTI-CLOUD & LLM

5.1 Customer Dashboard Insights

The customer dashboard grants users the functionality to sort properties and presents recommended properties obtained from Gemini language model suggestions. The customer sought a 2 BHK apartment for rent in Adyar that had to be exactly a 5-year-old South-East facing property along with gym amenities, security services, parking facilities and guaranteed water service.

The system advised selecting Property ID: RE29940196 following an assessment of both listings. The recommendation originated from a comparison between different features.

- The rental amount at ₹16,827 per month is less expensive than the potential alternative choice of ₹24,165 per month.
- Larger square footage (4700 sq.ft vs 3125 sq.ft)
- Superior amenities (security guards, gym, convention hall, bike parking)
- Better water access (corporation and borewell supply)
- Public sewer system over septic tank
- Central location in Adyar

A few minor drawbacks like laminate flooring and no balcony existed in the property yet they were balanced by multiple significant advantages. The LLM recommended both a 5–7% lower rent cost and a 10–15% smaller deposit amount to support customers in their purchasing process.

There are two benefits to customized property analysis done through AI which lead to both time efficiency and user confidence during selection processes. Users benefit from a summary that showcases essential trade-offs since AI technology both helps them make decisions through recommended actions.

User interactions with the customer dashboard enable a simplified process to search for properties. A customer who desires to rent a 2 BHK flat in Adyar with properties five years old should apply these search parameters from the table 5.1.

- Property Purpose: Rent
- Locality: Adyar

- House Type: Apartment
- BHK: 2
- Property Age: 5

Filtered Results

Based on this input, two properties matched the criteria:

ID	Rent	Sqft	Bathrooms	Floors	Amenities	Water	Parking	HOA
RE29940196	₹16,827	4700	3	10	Security Guards, Gym, Convention Hall	Both	Bike	₹25,425
RE67996354	₹24,165	3125	1	4	Police Station Nearby	Borewell	No	₹26,403

Table:5.1 Filtered Results

Gemini LLM Recommendation: RE29940196

Reasons:

- 50% larger area with 30% lower rent
- Comprehensive amenities and public sewer system
- Higher floor count indicates better facilities
- Both water supply and proper bike parking
- Greater bathroom availability and more supermarkets nearby

Concerns:

- Laminate flooring may wear out
- South-East facing may not appeal to all
- No balcony

Negotiation Advice:

- Negotiate monthly rent down 5–7% (~₹850–₹1,200)
- Try reducing the deposit by 10–15%
- Consider HOA fee absorption into monthly rent

Final Checks:

- Inspect flooring
- Review HOA coverage
- Verify water supply reliability

This customer-side recommendation system improves transparency and personalized property hunting, saving time and offering actionable insights from the fig 5.1 & 5.2.

The screenshot shows a web browser window with the URL `localhost:8501/customer`. The interface includes a sidebar with filters for 'Nearby Supermarkets' (set to 2), 'Property Style' (set to modern), and a search bar. Below the sidebar, a list of filtering conditions is displayed: 'Filtering: property_purpose == rent', 'Filtering: locality == adyar', 'Filtering: house_type == apartment', 'Filtering: bhk == 2', and 'Filtering: property_age == 5'. A 'Filtered Preview' section shows two rows of property data:

property_id	city	locality	property_purpose	house_type	bhk	house_floor	building_floor	property_age	facing	square_feet	rent_monthly_cost	owner_price_per_sqft	gov_price_per_sqft
11001	RE29940196	Chennai	Adyar	rent	apartment	2	7	10	5	South-East	4700	16827	0
15744	RE67996354	Chennai	Adyar	rent	apartment	2	3	4	5	North	3125	24165	0

A green banner at the bottom indicates 'Found 2 matching properties.'

property_id	city	locality	property_purpose	house_type	bhk	house_floor	building_floor	property_age	facing	square_feet	rent_monthly_cost	owner_price_per_sqft	gov_price_per_sqft
11001	RE29940196	Chennai	Adyar	rent	apartment	2	7	10	5	South-East	4700	16827	0
15744	RE67996354	Chennai	Adyar	rent	apartment	2	3	4	5	North	3125	24165	0

Fig:5.1 Customer Result page

The screenshot shows a web application interface for a real estate platform. At the top, there are several tabs including "realstate_project_report - Google Sheets", "Untitled spreadsheet - Google Sheets", "Real Estate Data Engineering", "New Tab", and "customer". The main content area has a header "Search Properties" with a "Deploy" button. Below it is a list of filtering criteria:

- Filtering: property_purpose == rent
- Filtering: locality == adyar
- Filtering: house_type == apartment
- Filtering: bhk == 2
- Filtering: property_age == 5

A "Filtered Preview" section shows a table with two rows of property data:

property_id	city	locality	property_purpose	house_type	bhk	house_floor	building_floor	property_age	facing	square_feet	rent_monthly_cost	owner_price_per_sqft	gov_price_per_sqft	security
11001	RE29940196	Chennai	Adyar	rent	apartment	2	7	10	5 South-East	4700	16827	0	0	Security Guards
15744	RE67996354	Chennai	Adyar	rent	apartment	2	3	4	5 North	3125	24165	0	0	Police Station Nearby

A green banner at the bottom left says "Found 2 matching properties." Below the table is an "OpenAI Recommendation" section:

Property Recommendation:
RE29940196

Reasoning:

1. RE29940196 offers a larger square footage (4700 sqft) compared to RE67996354 (3125 sqft) for a lower monthly rent (16827 vs 24165).
2. RE29940196 has more amenities including security guards, a gym, a convention hall, and proper parking (bike). RE67996354 has no parking available.
3. RE29940196 has a more comprehensive water supply (both) compared to RE67996354 (borewell only).
4. RE29940196 has more bathrooms available than RE67996354.

Fig:5.2 Customer LLm recommend response

5.2 Owner Dashboard Insights

The owner-side dashboard contains tools for predicting rent and buying prices of properties through property characteristic inputs. With this platform owners can add extensive details about their real estate properties encompassing type, measurements, levels, preparedness information and scheduled availability and nearby facilities. With these provided parameters two predictive models generated through Google Vertex AI training process using AWS S3 storage facilities provide estimated prices.

For example, in a test case for a buy scenario in Adyar:

- Predicted Price per Sqft: ₹13,760
- Total Estimated Property Price: ₹13,759,720

The estimate received additional evaluation through OpenAI's LLM which performs contextual analysis. The model determined that this property with its sea view and marble flooring as well as gym and swimming pool facilities was undervalued. The model proposed to increase the property cost by ₹1,240,000 by elevating the sqft price from ₹10,000 to ₹15,000. AI processing can generate essential

business intelligence that helps owners optimize their luxury market pricing by integrating market trends with customer feature requirements.

As part of the rental module owners need to enter comparable features which generates estimated monthly rental amounts. The rent model reached a predictive accuracy level stabilized at ₹3,189 RMSE while also demonstrating a reliable R² value of 0.818.

Through the owner dashboard users can obtain buy or rent price estimations for their properties by using comprehensive input fields alongside an ML processing system. Key input fields include:

- Location, house type, BHK, floor counts
- Age, facing direction, amenities (gym, pool, parking, security)
- Infrastructure like sewer, water, IT hub proximity
- Nearby schools, hospitals, supermarkets
- Interior features: flooring, roof, furnishing, photos
- Availability date and HOA fees

Example Input for Rent Estimation:

- Location: Adyar, Chennai
- Property: 2 BHK Apartment, 5 years old
- Amenities: Gym, Convention Hall, Security, Parking
- Water Supply: Both
- View: Sea
- HOA: ₹1,500
- Furnishing: Full

Model Output:

- Predicted Rent: ₹13,760/month
- LLM Suggests: Increase rent by ₹1,000–₹1,200/month for sea view and premium features

For Sale Prediction:

Predicted Price per Sqft: ₹13,760

Total Estimated Price: ₹13,759,720 (for 1,000 sqft)

LLM Advice:

- Increase to ₹15,000/sqft citing luxury amenities and location
- Expected profit: ₹1,240,000

This recommendation encourages owners to better price their property in line with market expectations and competitive offerings from the fig 5.3.

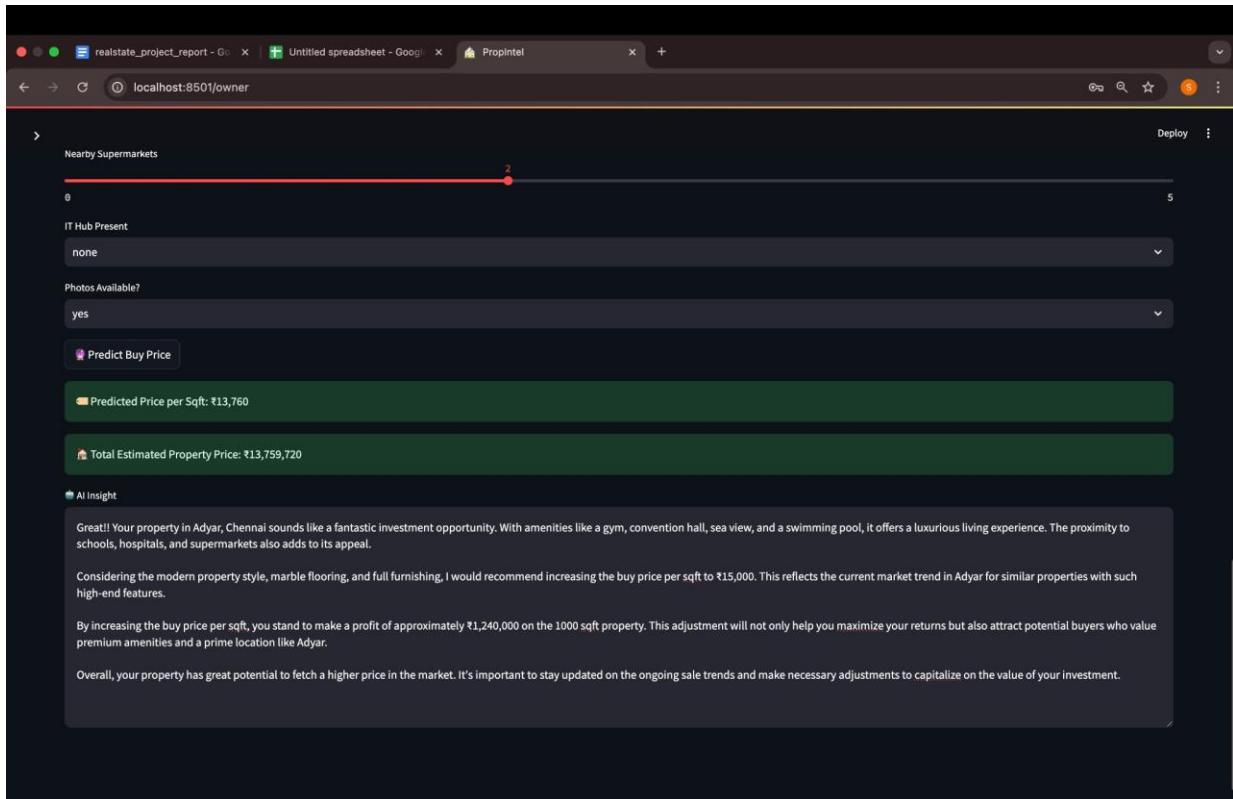


Fig:5.3 Owners model house prediction and LLm recommended response

5.3 Visual Insights with Amazon QuickSite

The transformed S3 bucket manages visualized datasets accessible to customers and property owners using Amazon QuickSight for data-driven decision-making. QuickSight demonstrates various types of interactive graphs through its interfaces that provide insights to users from the fig 5.4 & 5.5.

- Average rent and buy price trends across Chennai localities (e.g., Adyar, Anna Nagar, Velachery)
- Property value appreciation by locality and BHK configuration
- Distribution of properties by type (apartment, villa, etc.)

- Correlation of HOA fees with rent and buy prices
- Market availability calendar based on availability date inputs
- Heatmaps of amenities vs. property prices
- Filtering and drill-down capability based on city, locality, furnishing, and more from 5.9.

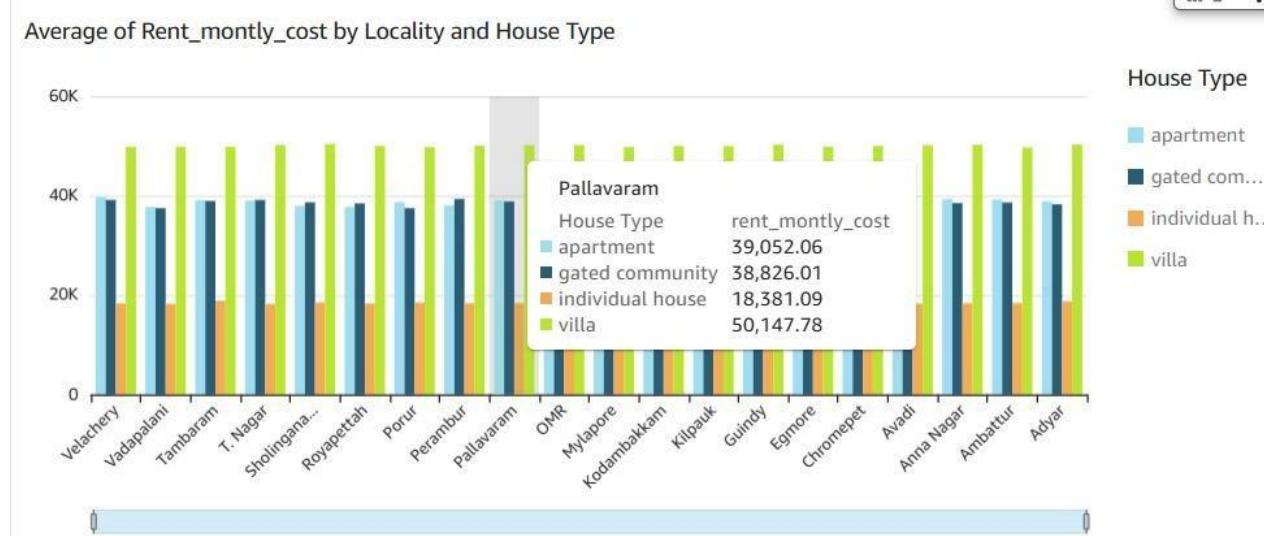


Fig:5.4 Average Rental Month Cost by Locality & House Type



Fig:5.5 Average Rent Month Cost by Locality & BHK

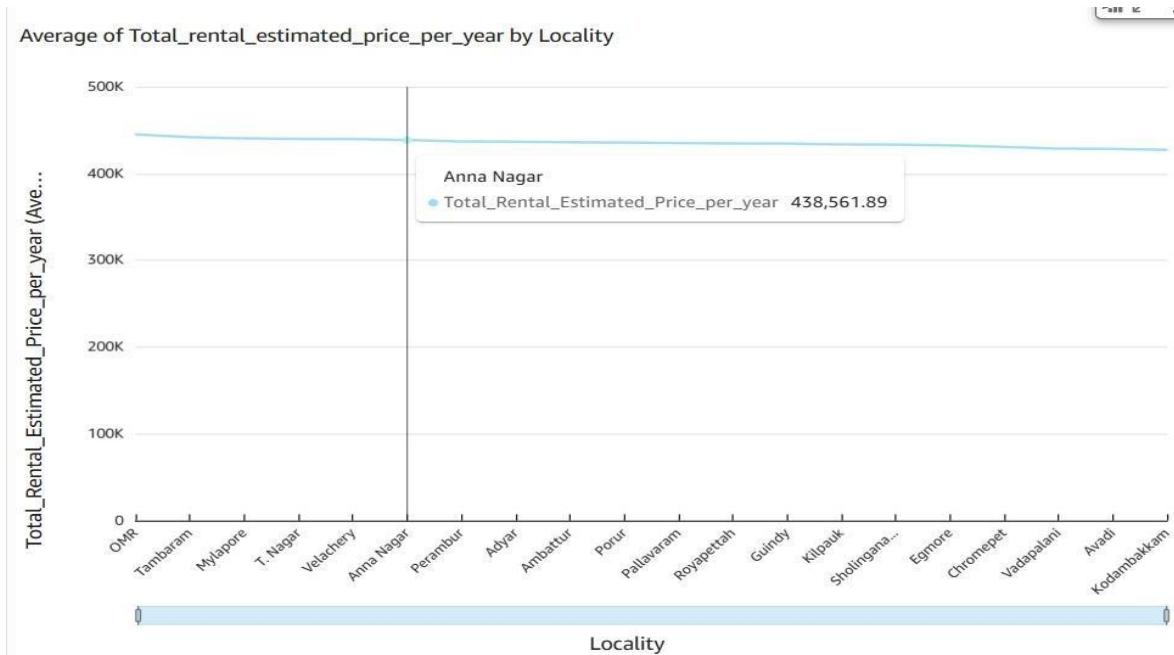


Fig:5.6 Total Rent Price per Year by Locality



Fig:5.7 Average Rent Cost by Availability Date,BHK & Locality

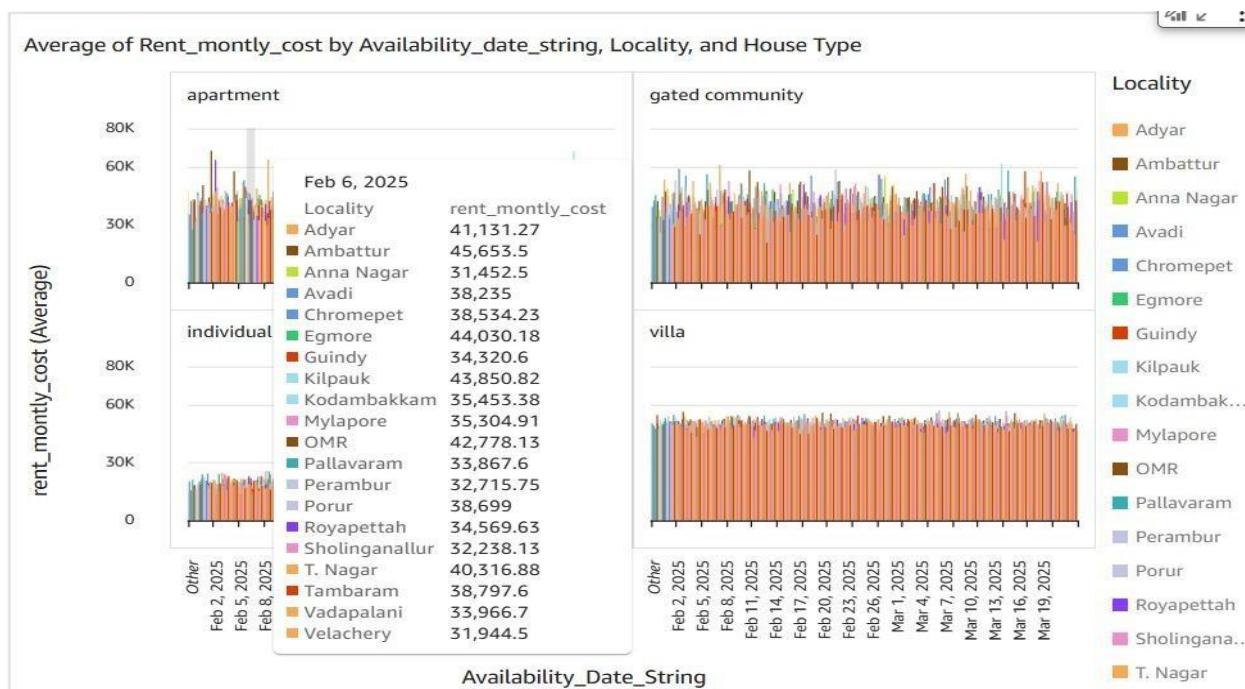


Fig:5.8 Average Rent Month Cost by Availability Date,Locality & House type

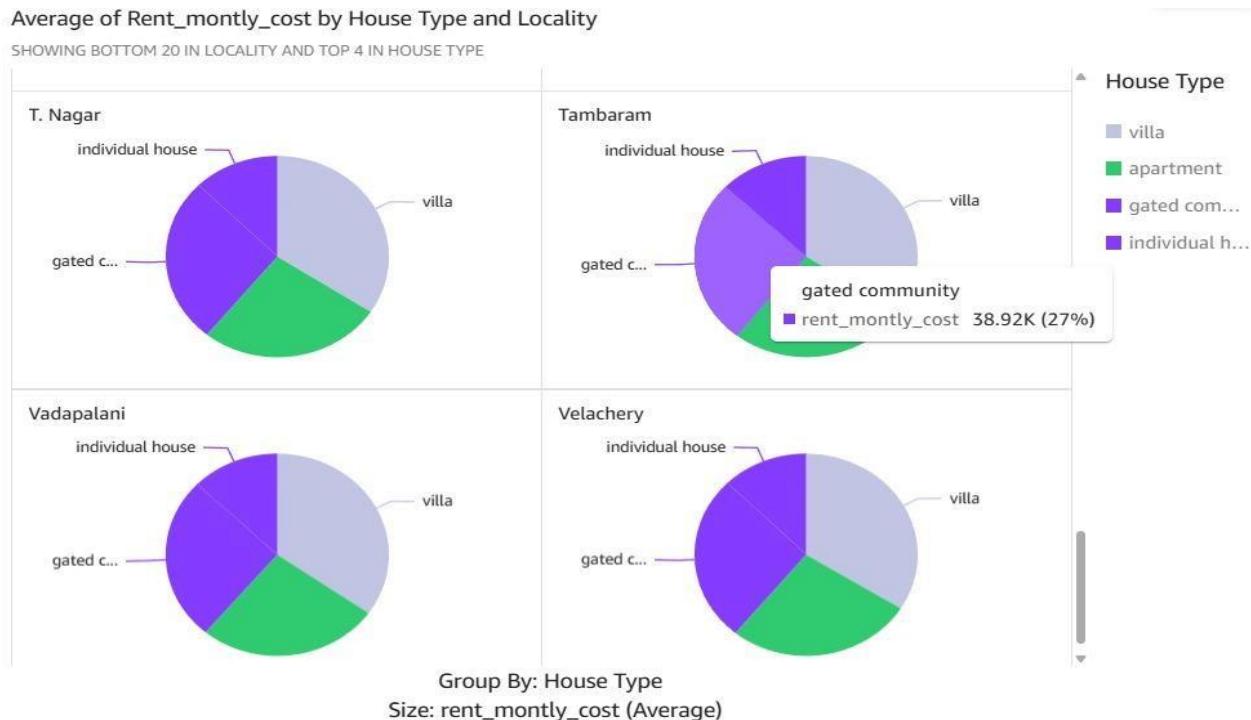


Fig:5.9 Average Rent Month Cost by House Type & Locality

CHAPTER 6

CONCLUSION

The developed real estate AI dashboard establishes effective connections between data-based analysis and user-centric system features targeting both property owners and customers in Chennai. The platform provides a user-friendly intelligent experience through the combination of model training from Google Vertex AI and deployment by AWS SageMaker and visualization through Amazon QuickSight for real-world property requirements.

Users can access precise rent and purchase cost forecasts alongside LLM-based recommendation services which provide practical negotiation guidelines through the customer dashboard. Through price recommendations generated from data the system enables owners to see both present market trends which helps them determine competitive rental rates.

A reliable prediction framework exists because the machine learning model achieves an R^2 value of 0.818 and an RMSE of ₹3,188 while QuickSight visuals enable users to simplify complex data analysis. The system provides complete knowledge through its integration of amenities alongside locality trends and infrastructure quality and property view and facing assessments.

Rich data processing through intelligent design systems controlled by artificial intelligence technology enhances the real estate industry by generating reliable stakeholder decisions with minimized obstacles.

CHAPTER 7

FUTURE WORK

The current real estate dashboard stands well as a base for AI-based property intelligence but further adjustments will enhance both estimation precision and network functionality and application capabilities.

Modern geospatial data can be integrated through mapping APIs (for example Google Maps and Mapbox) which improve location-based analysis. Users benefit from enhanced purchasing decisions through the visual analysis of location distances toward schools and hospitals and public transportation and commercial centers. Commute time calculation implementation will benefit all stakeholders in the platform.

Maps through API integration like Google Maps and Mapbox let users combine contemporary geospatial information to enhance their location-based assessment capabilities. Users obtain better purchasing options by using visualizations for gauging proximity between schools and hospitals and public transport stations and commercial outlets. The implementation of commute time calculation services will produce advantages for all platform stakeholders.

A feedback system that collects user transaction prices for the system will enable both improved model accuracy along with personalized future recommendations. The accessibility for non-English speakers can be enhanced through both language extension and development of voice-based search capabilities.

A feedback system that collects user transaction prices for the system will enable both improved model accuracy along with personalized future recommendations. The accessibility for non-English speakers can be enhanced through both language extension and development of voice-based search capabilities.

REFERENCES

- [1] N. N. Ghosalkar and S. N. Dhage, "Real Estate Value Prediction Using Linear Regression," 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBE), Pune, India, 2018, pp. 1-5, doi: 10.1109/ICCUBE.2018.8697639.
- [2] W. Aiyan and X. Yanmei, "Multiple Linear Regression Analysis of Real Estate Price," 2018 International Conference on Robots & Intelligent System (ICRIS), Changsha, China, 2018, pp. 564-568, doi: 10.1109/ICRIS.2018.00145.
- [3] C. R. Madhuri, G. Anuradha and M. V. Pujitha, "House Price Prediction Using Regression Techniques: A Comparative Study," 2019 International Conference on Smart Structures and Systems (ICSSS), Chennai, India, 2019, pp. 1-5, doi: 10.1109/ICSSS.2019.8882834.
- [4] A. Mukherjee, R. Nambiar and D. Choudhury, "Machine Learning based Real Estate Price Prediction," 2024 8th International Conference on Inventive Systems and Control (ICISC), Coimbatore, India, 2024, pp. 351-358, doi: 10.1109/ICISC62624.2024.00067.
- [5] A. Varma, A. Sarma, S. Doshi and R. Nair, "House Price Prediction Using Machine Learning and Neural Networks," 2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT), Coimbatore, India, 2018, pp. 1936-1939, doi: 10.1109/ICICCT.2018.8473231.
- [6] A. Kumar, K. Singh, A. Raj and A. Shakya, "House Price Prediction Using ML," 2024 International Conference on Sustainable Power & Energy (ICSPE), Raigarh, India, 2024, pp. 1-5, doi: 10.1109/ICSPE62629.2024.10924379.
- [7] Q. Chen, H. Zhang and B. Zhao, "Research on Real Estate Price Analysis and Prediction Based on Genetic Algorithm," 2024 Asia-Pacific Conference on Software Engineering, Social Network Analysis and Intelligent Computing (SSAIC), New Delhi, India, 2024, pp. 627-630, doi: 10.1109/SSAIC61213.2024.00126.
- [8] L. Bo and Z. Ming-yu, "Mathematical Analysis of Relationship between Land Price and Real Estate Bubbles," 2006 International Conference on Management Science and Engineering, Lille, France, 2006, pp. 2167-2171, doi: 10.1109/ICMSE.2006.314151.
- [9] B. Singh, J. Singh, A. Kumar and S. Gupta, "Investigating the Effectiveness of Machine Learning Techniques for Real Estate Price Prediction," 2023 4th International Conference on Smart Electronics and Communication (ICOSEC), Trichy, India, 2023, pp. 1684-1688, doi: 10.1109/ICOSEC58147.2023.10275956.

- [10] B. Singh, J. Singh, A. Kumar and S. Gupta, "Investigating the Effectiveness of Machine Learning Techniques for Real Estate Price Prediction," 2023 4th International Conference on Smart Electronics and Communication (ICOSEC), Trichy, India, 2023, pp. 1684-1688, doi: 10.1109/ICOSEC58147.2023.10275956.
- [11] Kamath, U., Keenan, K., Somers, G., & Sorenson, S. (2024). LLMs in Production. In Large Language Models: A Deep Dive: Bridging Theory and Practice (pp. 315-373). Cham: Springer Nature Switzerland.
- [12] Nguyen, K. H., Issa, F., Markar, A., Ahmed, S., & Mahato, B. Develop Machine Learning Models using AWS Sagemaker.
- [13] Nigenda, D., Karnin, Z., Zafar, M. B., Ramesha, R., Tan, A., Donini, M., & Kenthapadi, K. (2022, August). Amazon sagemaker model monitor: A system for real-time insights into deployed machine learning models. In Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (pp. 3671-3681).
- [14] Chen, T., & Si, S. (2024). Predicting Rental Price of Lane Houses in Shanghai with Machine Learning Methods and Large Language Models. arXiv preprint arXiv:2405.17505.
- [15] Ma, Y., Zhang, T., & Zhan, G. (2024, June). An LLM-based Intelligent System for the Evaluation of Property Geographical Environment. In 2024 International Symposium on Intelligent Robotics and Systems (ISoIRS) (pp. 258-262). IEEE.
- [16] Wan, H., Zhang, J., Chen, Y., Xu, W., & Feng, F. (2024). Generative AI Application for Building Industry. arXiv preprint arXiv:2410.01098.
- [17] Sandeep Kumar, E., Talasila, V., Rishe, N., Suresh Kumar, T. V., & Iyengar, S. S. (2019). Location identification for real estate investment using data analytics. International Journal of Data Science and Analytics, 8(3), 299-323.
- [18] Zhao, Y., & Gao, H. (2024). Utilizing large language models for information extraction from real estate transactions. arXiv preprint arXiv:2404.18043.
- [19] Kvet, M., Potoèár, M., & Tatarka, S. (2025). Real Estate Attribute Value Extraction Using Large Language Models. IEEE Access.
- [20] Naushad, R., Gupta, R., Bhutiyal, T., & Prajapati, V. (2024, May). A Novel Approach to Rental Market Analysis for Property Management Firms Using Large Language Models and Machine Learning. In International Joint Conference on Rough Sets (pp. 247-261). Cham: Springer Nature Switzerland.

PROJECT FILES AND DEMO REFERENCES

Github:

<https://github.com/Jsjanchowdary/REAL-ESTATE-DATA-ENGINEERING-USING-MULTI-CLOUD---LLM-project>

Youtube Demo: <https://youtu.be/JA1kmpehzBM>

APPENDIX

Model file code:

```
pip install pyspark
pip install boto3
import os
import boto3
os.environ["AWS_ACCESS_KEY_ID"] = "*****"
os.environ["AWS_SECRET_ACCESS_KEY"] = "*****"
s3 = boto3.client('s3')
bucket = "magedataset"
key = "real_state_data.csv"
local_file = "/tmp/real_state_data.csv"
# Download the file
s3.download_file(bucket, key, local_file)
print("Download complete.")
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName("RealEstateData").getOrCreate()
df = spark.read.csv(local_file, header=True, inferSchema=True)
df.printSchema()
df.show(5)
Property_ID: string (nullable = true)
|-- City: string (nullable = true)
|-- Locality: string (nullable = true)
|-- Property Purpose: string (nullable = true)
|-- House Type: string (nullable = true)
|-- BHK: integer (nullable = true)
|-- House Floor: integer (nullable = true)
|-- Building Floor: integer (nullable = true)
|-- Property Age: integer (nullable = true)
|-- Facing: string (nullable = true)
|-- Square Feet: integer (nullable = true)
|-- rent_montly_cost: integer (nullable = true)
|-- Owner Price per Sqft: integer (nullable = true)
|-- GOV Price per Sqft: integer (nullable = true)
|-- Security: string (nullable = true)
|-- Gym: string (nullable = true)
|-- Convention Hall: string (nullable = true)
|-- Parking: string (nullable = true)
|-- Water Supply: string (nullable = true)
|-- Bathrooms: integer (nullable = true)
|-- Balcony: integer (nullable = true)
|-- Availability Date: date (nullable = true)
|-- Deposit Cost: integer (nullable = true)
|-- Furnishing: string (nullable = true)
|-- HOA Fees : double (nullable = true)
|-- Flooring Type: string (nullable = true)
```

```

|-- Roof Type: string (nullable = true)
|-- Property Style: string (nullable = true)
|-- View: string (nullable = true)
|-- Nearby Schools: integer (nullable = true)
|-- Nearby Hospitals: integer (nullable = true)
|-- Nearby Supermarkets: integer (nullable = true)
|-- IT Hub Present: string (nullable = true)
|-- Swimming Pool: string (nullable = true)
|-- Lot Shape: string (nullable = true)
|-- Sewer System: string (nullable = true)
|-- Address: string (nullable = true)
|-- owner_name: string (nullable = true)
|-- owner_email: string (nullable = true)
|-- Owner Phone: string (nullable = true)
|-- Pincode: integer (nullable = true)
|-- State: string (nullable = true)
|-- Photos: string (nullable = true)
from pyspark.sql.functions import trim, lower, col
string_columns = [col_name for col_name, dtype in df.dtypes if dtype == 'string']
for col_name in string_columns:
    df = df.withColumn(col_name, trim(lower(col(col_name))))
# Replace empty strings with nulls
from pyspark.sql.functions import when
for c in df.columns:
    df = df.withColumn(c, when(col(c) == "", None).otherwise(col(c)))
# Fill missing numeric values with 0 or median (simplified here as 0)
num_cols = [col_name for col_name, dtype in df.dtypes if dtype in ['int', 'double']]
df = df.fillna(0, subset=num_cols)
# Fill string columns with "unknown"
df = df.fillna("unknown", subset=string_columns)
# Step 3: Partition
buyer_df = df.filter(col("Property Purpose") == "buy").drop("rent_monthly_cost", "Deposit Cost")
rental_df = df.filter(col("Property Purpose") == "rent").drop("Owner Price per Sqft", "GOV Price per Sqft")
from pyspark.ml.feature import StringIndexer, OneHotEncoder, VectorAssembler
from pyspark.ml import Pipeline
# Define columns
numeric_features = [
    "BHK", "House Floor", "Building Floor", "Property Age", "Square Feet",
    "Bathrooms", "Balcony", "HOA Fees ", "Nearby Schools", "Nearby Hospitals",
    "Nearby Supermarkets"
]
categorical_features = [
    "City", "Locality", "Property Purpose", "House Type", "Facing", "Security", "Gym",
    "Convention Hall", "Parking", "Water Supply", "Furnishing", "Flooring Type",
    "Roof Type", "Property Style", "View", "IT Hub Present", "Swimming Pool",
    "Lot Shape", "Sewer System"
]
# You may drop or convert Availability Date and Photos later if needed
from pyspark.sql.functions import col, months_between, current_date, when, lower

```

```

# 1. BHK × Sqft
ml_df = buyer_df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))

# 2. Months to Availability
ml_df = ml_df.withColumn("Months_To_Availability", months_between(col("Availability Date"), current_date()))

# 3. Has Photos
ml_df = ml_df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))

# 4. Amenity Score (gym, pool, security, convention hall)
ml_df = ml_df.withColumn("Amenity_Score",
    (when(lower(col("Gym")) == "yes", 1).otherwise(0) +
     when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
     when(lower(col("Security")) != "none", 1).otherwise(0) +
     when(lower(col("Convention Hall")) == "yes", 1).otherwise(0))
)

# 5. Is Near IT Hub
ml_df = ml_df.withColumn("Is_Near_IT_Hub", when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))

new_numeric_features = [
    "BHK_Sqft", "Months_To_Availability", "Amenity_Score", "Has_Photos", "Is_Near_IT_Hub"
]
numeric_features += new_numeric_features

# Index and OneHotEncode categorical features
indexers = [StringIndexer(inputCol=col, outputCol=col + "_Index", handleInvalid='keep') for col in categorical_features]
encoders = [OneHotEncoder(inputCol=col + "_Index", outputCol=col + "_OHE") for col in categorical_features]

# Final input features (numeric + encoded)
encoded_features = [col + "_OHE" for col in categorical_features]
assembler_inputs = numeric_features + encoded_features
assembler = VectorAssembler(
    inputCols=assembler_inputs,
    outputCol="features"
)

# Target column
target_column = "Owner Price per Sqft"
pipeline = Pipeline(stages=indexers + encoders + [assembler])
pipeline_model = pipeline.fit(ml_df)
processed_df = pipeline_model.transform(ml_df)
train_df, test_df = processed_df.randomSplit([0.8, 0.2], seed=42)
from pyspark.ml.regression import LinearRegression
lr = LinearRegression(
    featuresCol="features",
    labelCol=target_column,
    predictionCol="prediction"
)
lr_model = lr.fit(train_df)
results = lr_model.evaluate(test_df)
print("RMSE:", results.rootMeanSquaredError)
print("R2:", results.r2)
RMSE: 3188.8767498553666
R2: 0.8178149248860942
user_input = {

```

```

    "City": "chennai",
    "Locality": "adyar",
    "Property Purpose": "buy",
    "House Type": "villa",
    "BHK": 4,
    "House Floor": 0,
    "Building Floor": 2,
    "Property Age": 3,
    "Facing": "north-east",
    "Square Feet": 3200,
    "Security": "cctv",
    "Gym": "no",
    "Convention Hall": "no",
    "Parking": "car",
    "Water Supply": "corporation",
    "Bathrooms": 3,
    "Balcony": 3,
    "Availability Date": "2025-04-10",
    "Furnishing": "full",
    "HOA Fees": 1500.0,
    "Flooring Type": "tiles",
    "Roof Type": "concrete",
    "Property Style": "modern",
    "View": "city view",
    "Nearby Schools": 4,
    "Nearby Hospitals": 4,
    "Nearby Supermarkets": 5,
    "IT Hub Present": "tcs",
    "Swimming Pool": "yes",
    "Lot Shape": "rectangular",
    "Sewer System": "public",
    "Photos": "yes"
}
from pyspark.sql import Row
from pyspark.sql.functions import col, when, months_between, current_date, lower
# 1. Create the user input DataFrame
user_row = Row(**user_input)
user_df = spark.createDataFrame([user_row])
# 2. Manually add engineered features
user_df = user_df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
user_df = user_df.withColumn("Months_To_Availability",
                             months_between(col("Availability Date"), current_date()))
user_df = user_df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
user_df = user_df.withColumn("Amenity_Score",
                             (when(lower(col("Gym")) == "yes", 1).otherwise(0) +
                              when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
                              when(lower(col("Security")) != "none", 1).otherwise(0) +
                              when(lower(col("Convention Hall")) == "yes", 1).otherwise(0)))
)

```

```

user_df = user_df.withColumn("Is_Near_IT_Hub",
                             when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
# Transform and predict
user_transformed = pipeline_model.transform(user_df)
prediction = lr_model.transform(user_transformed)
prediction.select("prediction").show()
-----
|   prediction|
+-----+
|20715.184733821|
+-----+
from pyspark.sql.functions import round as spark_round
# 2. Add Total_Predicted_Price column
prediction = prediction.withColumn(
    "Total_Predicted_Price",
    spark_round(col("prediction") * col("Square Feet")))
)
# 3. Get values from prediction DataFrame
row = prediction.select("prediction", "Square Feet", "Total_Predicted_Price").first()
# 4. Extract values
predicted_price_per_sqft = round(row["prediction"])
total_price = row["Total_Predicted_Price"]
# □ Display nicely
print(f"□ Predicted Price per Sqft: ₹{predicted_price_per_sqft:,}")
print(f"□ Total Estimated Property Price: ₹{total_price:,}")
Predicted Price per Sqft: ₹20,715
□ Total Estimated Property Price: ₹66,288,591.0
from pyspark.sql.functions import round as spark_round
# 2. Add Total_Predicted_Price column
prediction = prediction.withColumn(
    "Total_Predicted_Price",
    spark_round(col("prediction") * col("Square Feet")))
)
# 3. Get values from prediction DataFrame
row = prediction.select("prediction", "Square Feet", "Total_Predicted_Price").first()
# 4. Extract values
predicted_price_per_sqft = round(row["prediction"])
total_price = row["Total_Predicted_Price"]
# □ Display nicely
print(f"□ Predicted Price per Sqft: ₹{predicted_price_per_sqft:,}")
print(f"□ Total Estimated Property Price: ₹{total_price:,}")
Predicted Price per Sqft: ₹20,768
□ Total Estimated Property Price: ₹66,457,444.0 #
Save regression model
lr_model.save("house_price_model")
# Save preprocessing pipeline
pipeline_model.save("preprocessing_pipeline")
from pyspark.sql.functions import col
rental_df = rental_df.filter(col("Property Purpose") == "rent")

```

```

numeric_features = [
    "BHK", "House Floor", "Building Floor", "Property Age", "Square Feet",
    "Bathrooms", "Balcony", "Deposit Cost", "HOA Fees ",
    "Nearby Schools", "Nearby Hospitals", "Nearby Supermarkets"
]
categorical_features = [
    "City", "Locality", "House Type", "Facing", "Security", "Gym", "Convention Hall",
    "Parking", "Water Supply", "Furnishing", "Flooring Type", "Roof Type",
    "Property Style", "View", "IT Hub Present", "Swimming Pool",
    "Lot Shape", "Sewer System"
]
from pyspark.sql.functions import months_between, current_date, when, lower
rental_df = rental_df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
rental_df = rental_df.withColumn("Months_To_Availability",
    months_between(col("Availability Date"), current_date()))
rental_df = rental_df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
rental_df = rental_df.withColumn("Amenity_Score",
    (when(lower(col("Gym")) == "yes", 1).otherwise(0) +
     when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
     when(lower(col("Security")) != "none", 1).otherwise(0) +
     when(lower(col("Convention Hall")) == "yes", 1).otherwise(0)))
)
rental_df = rental_df.withColumn("Is_Near_IT_Hub",
    when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
numeric_features += ["BHK_Sqft", "Months_To_Availability", "Has_Photos", "Amenity_Score",
    "Is_Near_IT_Hub"]
from pyspark.ml.feature import StringIndexer, OneHotEncoder, VectorAssembler
from pyspark.ml import Pipeline
# Step 4.1: Index and Encode Categorical
indexers = [StringIndexer(inputCol=c, outputCol=f"{c}_Index", handleInvalid="keep") for c in categorical_features]
encoders = [OneHotEncoder(inputCol=f"{c}_Index", outputCol=f"{c}_OHE") for c in categorical_features]
# Step 4.2: Assemble Features
encoded_features = [f"{c}_OHE" for c in categorical_features]
assembler_inputs = numeric_features + encoded_features
assembler = VectorAssembler(inputCols=assembler_inputs, outputCol="features")
# Step 4.3: Pipeline
pipeline = Pipeline(stages=indexers + encoders + [assembler])
pipeline_model = pipeline.fit(rental_df)
processed_df = pipeline_model.transform(rental_df)
train_df, test_df = processed_df.randomSplit([0.8, 0.2], seed=42)
from pyspark.ml.regression import LinearRegression
lr = LinearRegression(
    featuresCol="features",
    labelCol="rent_montly_cost",
    predictionCol="prediction"
)
lr_model = lr.fit(train_df)
results = lr_model.evaluate(test_df)

```

```

print("❖ Rental Model Performance")
print("RMSE:", results.rootMeanSquaredError)
print("R²:", results.r2)
❖Rental Model Performance
RMSE: 11004.180780888408
R²: 0.5816076317772093
user_input_rent = {
    "City": "chennai",
    "Locality": "velachery",
    "Property Purpose": "rent",
    "House Type": "apartment",
    "BHK": 2,
    "House Floor": 3,
    "Building Floor": 5,
    "Property Age": 8,
    "Facing": "north-east",
    "Square Feet": 1700,
    "Security": "cctv",
    "Gym": "yes",
    "Convention Hall": "no",
    "Parking": "both",
    "Water Supply": "corporation",
    "Bathrooms": 2,
    "Balcony": 1,
    "Availability Date": "2025-04-15",
    "Deposit Cost": 50000,
    "Furnishing": "semi",
    "HOA Fees": 1800.0,
    "Flooring Type": "tiles",
    "Roof Type": "concrete",
    "Property Style": "modern",
    "View": "city view",
    "Nearby Schools": 3,
    "Nearby Hospitals": 2,
    "Nearby Supermarkets": 3,
    "IT Hub Present": "tcs",
    "Swimming Pool": "no",
    "Lot Shape": "rectangular",
    "Sewer System": "public",
    "Photos": "yes"
}
from pyspark.sql import Row
from pyspark.sql.functions import col, months_between, current_date, when, lower
# Convert to DataFrame
user_row = Row(**user_input_rent)
user_rent_df = spark.createDataFrame([user_row])
# Add engineered features
user_rent_df = user_rent_df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
user_rent_df = user_rent_df.withColumn("Months_To_Availability",

```

```

months_between(col("Availability Date"), current_date()))
user_rent_df = user_rent_df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
user_rent_df = user_rent_df.withColumn("Amenity_Score",
    (when(lower(col("Gym")) == "yes", 1).otherwise(0) +
     when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
     when(lower(col("Security")) != "none", 1).otherwise(0) +
     when(lower(col("Convention Hall")) == "yes", 1).otherwise(0)))
)
user_rent_df = user_rent_df.withColumn("Is_Near_IT_Hub",
    when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
# Transform input
user_transformed = pipeline_model.transform(user_rent_df)
# Predict
prediction = lr_model.transform(user_transformed)
# Format Output
row = prediction.select("prediction").first()
predicted_rent = round(row["prediction"])
print(f"\u25aa Predicted Monthly Rent: ₹{predicted_rent:,}")
Predicted Monthly Rent: ₹33,203
!zip -r house_price_model.zip house_price_model
!zip -r preprocessing_pipeline.zip preprocessing_pipeline
!zip -r rental_model.zip rental_model
!zip -r rental_pipeline.zip rental_pipeline

```

Main code:

```

import streamlit as st
import pandas as pd
import hashlib
import os
import requests
import streamlit_lottie
# Set the page to default to wide mode
st.set_page_config(page_title="PropIntel", page_icon="\u25aa", layout="wide")
def load_lottieur(url):
    r = requests.get(url)
    if r.status_code != 200:
        return None
    return r.json()
# CSV file paths
CUSTOMER_FILE = "/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/customer_data.csv"
ADMIN_FILE = "/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/admin_owner_data.csv"
def main():
    # Initialize session state
    if "logged_in" not in st.session_state:
        st.session_state["logged_in"] = False
        st.session_state["user_role"] = None

```

```

st.session_state["user_name"] = None
st.session_state["page"] = "Login" # Default page
def initialize_csv(file_path):
    """Ensures CSV file exists and has proper headers."""
    if not os.path.exists(file_path) or os.stat(file_path).st_size == 0:
        df = pd.DataFrame(columns=["name", "email", "password"])
        df.to_csv(file_path, index=False)
initialize_csv(CUSTOMER_FILE)
initialize_csv(ADMIN_FILE)
def hash_password(password):
    return hashlib.sha256(password.encode()).hexdigest()
def load_user_data(file_path):
    """Loads user data from CSV and ensures proper headers exist."""
    if os.path.exists(file_path) and os.stat(file_path).st_size > 0:
        return pd.read_csv(file_path)
    else:
        return pd.DataFrame(columns=["name", "email", "password"]) # ✅Safe empty DataFrame
def check_login(file_path, identifier, password):
    df = load_user_data(file_path)
    hashed_password = hash_password(password)
    if identifier in df["email"].values or identifier in df["name"].values:
        user_row = df[(df["email"] == identifier) | (df["name"] == identifier)]
        if not user_row.empty and user_row["password"].values[0] == hashed_password:
            return True, user_row["name"].values[0]
    return False, None
# ② If logged in, switch to the respective page if
st.session_state["logged_in"]:
    if st.session_state["user_role"] == "Customer":
        st.switch_page("/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/pages/customer.py") # ✅
    Redirect to Customer Dashboard
    elif st.session_state["user_role"] == "Admin Owner":
        st.switch_page("/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/pages/admin.py") # ✅
    Redirect to Admin Dashboard
    # ② Hide Sidebar Before Login
    if not st.session_state["logged_in"]:
        st.markdown("""
<style>
[data-testid="stSidebar"] { visibility: hidden; }
</style>
""", unsafe_allow_html=True)

# Center align using HTML
st.markdown(
"""
<h1 style='text-align: center;'>PropIntel</h1>
<h3 style='text-align: center;'>Smart Insights, Smarter Investments.</h3>
""",
unsafe_allow_html=True
)

```

```

st.title(" ")
col1, col2 = st.columns(2)
with col1:
    l1 = "https://lottie.host/5ae01eac-69f3-4f8b-8703-137ca5bbfc31/cSY7ipxxfn.json"
    st.lottie(l1)
with col2:
    # Display login or signup form based on session state
    if st.session_state["page"] == "Login":
        st.subheader("Login to Your Account")
        identifier = st.text_input("Enter your name or email")
        password = st.text_input("Enter your password", type="password")
        role = st.radio("Select Role", ["Customer", "Admin Owner"], horizontal=True)

    col3, col4 = st.columns(2)
    with col3:
        if st.button("Login", use_container_width=True):
            file_path = CUSTOMER_FILE if role == "Customer" else ADMIN_FILE
            login_success, user_name = check_login(file_path, identifier, password)
            if login_success:
                st.session_state["logged_in"] = True
                st.session_state["user_role"] = role # ↗ Set user role correctly
                st.session_state["user_name"] = user_name
                st.success(f"Login successful! Redirecting to {role} page...")
                # ↗ Redirect after successful login
            if role == "Customer":
                st.switch_page("/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/pages/customer.py")
            elif role == "Admin Owner":
                st.switch_page("/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/pages/admin.py")
            else:
                st.error("Invalid credentials! Please check your details.")
    with col4:
        # Button to switch to Sign-Up form
        if st.button("Create a New Account", use_container_width=True):
            st.session_state["page"] = "Sign Up"
            st.rerun()
    elif st.session_state["page"] == "Sign Up":
        st.subheader("Create a New Account (Customers Only)")
        name = st.text_input("Enter your name")
        email = st.text_input("Enter your email")
        password = st.text_input("Enter your password", type="password")
        re_password = st.text_input("Re-enter your password", type="password")
    col5, col6 = st.columns(2)
    with col5:
        if st.button("Sign Up", use_container_width=True):
            if password != re_password:
                st.error("Passwords do not match!")
            else:
                df = load_user_data(CUSTOMER_FILE)
                if email in df["email"].values:

```

```

        st.warning("This email is already registered!")
    else:
        new_user = pd.DataFrame([[name, email, hash_password(password)]],
                               columns=["name", "email", "password"])
        df = pd.concat([df, new_user], ignore_index=True)
        df.to_csv(CUSTOMER_FILE, index=False)
        st.success("Sign up successful! You can now log in.")
        # Switch back to login page
        st.session_state["page"] = "Login"
        st.rerun()
with col6:
    # Button to switch back to Login form
    if st.button("Back to Login", use_container_width=True):
        st.session_state["page"] = "Login"
        st.rerun()
if __name__ == "__main__":
    main()

```

Admin Code:

```

import streamlit as st
import requests
import streamlit_lottie
import pandas as pd
import boto3
import json
from io import StringIO
import uuid
st.set_page_config(page_title="PropIntel", page_icon="⚠️", layout="wide")
def load_lottieurl(url):
    r = requests.get(url)
    if r.status_code != 200:
        return None
    return r.json()
l1 = "https://lottie.host/20ddb2ba-92af-4a75-8efd-8e04beeb5dd7/vZ9Xa4h9ug.json"
# ⚠️ Hide Sidebar for Admin Dashboard
st.markdown("""
<style>
[data-testid="stSidebar"] { visibility: hidden; }
.logout-container {
    position: absolute;
    top: 10px;
    right: 20px;
    z-index: 1000;
}
</style>
""", unsafe_allow_html=True)
# ⚠️ Ensure User is Logged In
if "logged_in" not in st.session_state or not st.session_state["logged_in"]:

```

```

st.warning("You must be logged in to access this page.")
if st.button("Go to Login", use_container_width=True):
    st.switch_page("/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/main.py") # ↗ Redirect to Login
Page
    st.stop()
# st.success(f"Welcome, {st.session_state['user_name']}! You are logged in as an Admin.")
# ☒ Move Logout Button to the Top Right
logout_placeholder = st.empty()
with logout_placeholder.container():
    col1, col2 = st.columns([7, 1]) # Right-align using columns
    with col1:
        st.title(f"Hello, {st.session_state['user_name']}!")
    with col2:
        if st.button("Logout", use_container_width=True):
            st.session_state.clear() # ↗ Clears all session data
            st.switch_page("/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/main.py") # ↗ Redirect to
login page
    col3, col4 = st.columns(2)
    with col4:
        st.lottie(l1)
# AWS S3 Configuration
S3_BUCKET = "magedataset"
S3_FILE_KEY = "real_state_data.csv"
AWS_ACCESS_KEY = "keys"
AWS_SECRET_KEY = "keys"
AWS_REGION = "us-east-1"
# Function to Load Data from S3
def load_data_from_s3():
    s3 = boto3.client(
        "s3",
        aws_access_key_id=AWS_ACCESS_KEY,
        aws_secret_access_key=AWS_SECRET_KEY,
        region_name=AWS_REGION,
    )
    try:
        obj = s3.get_object(Bucket=S3_BUCKET, Key=S3_FILE_KEY)
        df = pd.read_csv(obj["Body"])
        return df
    except Exception as e:
        st.error(f"Error loading data: {e}")
        return pd.DataFrame()
# Function to Save Data to S3
def save_data_to_s3(df):
    s3 = boto3.client(
        "s3",
        aws_access_key_id=AWS_ACCESS_KEY,
        aws_secret_access_key=AWS_SECRET_KEY,
        region_name=AWS_REGION,
    )

```

```

csv_buffer = StringIO()
df.to_csv(csv_buffer, index=False)
s3.put_object(Bucket=S3_BUCKET, Key=S3_FILE_KEY, Body=csv_buffer.getvalue())
st.success("Property data saved successfully!")

# Function to Generate Unique Property ID
def generate_unique_property_id(existing_df):
    while True:
        new_id = f"PROP-{uuid.uuid4().hex[:6].upper()}" # Generates a random 6-character alphanumeric ID
        if new_id not in existing_df["Property_ID"].values:
            return new_id

# Streamlit UI
st.title("Owner Property Data Entry")
# Load existing data
df = load_data_from_s3()
with st.expander("Basic Property Details", expanded=True):
    city = st.selectbox("City", ["Chennai"], index=0)
    locality = st.selectbox("Locality", ["Adyar", "Anna Nagar", "T. Nagar", "Velachery", "OMR", "Tambaram", "Porur", "Guindy", "Mylapore", "Kilpauk", "Perambur", "Vadapalani", "Kodambakkam", "Royapettah", "Egmore", "Pallavaram", "Avadi", "Chromepet", "Sholinganallur", "Ambattur"], index=0)
    property_purpose = st.radio("Property Purpose", ["Rent", "Sale"], horizontal=True)
    house_type = st.selectbox("House Type", ["apartment", "individual house", "villa", "gated community"], index=0)
    bhk = st.selectbox("BHK", list(range(1, 7)), index=2)
    sqft = st.number_input("Square Feet", min_value=100, step=10, value=1000)
    sqft_cost = st.number_input("Cost per Square Feet", min_value=1000, step=100, value=5000)
with st.expander("Building Details"):
    house_floor = st.selectbox("House Floor", list(range(1, 16)), index=0)
    building_floor = st.selectbox("Building Floor", list(range(1, 16)), index=4)
    property_age = st.number_input("Property Age (years)", min_value=0, step=1, value=5)
    facing = st.selectbox("Facing Direction", ["North", "South", "East", "West", "North-East", "North-West", "South-East", "South-West"], index=0)
    photos = st.text_input("past the image links")
    availability_date = st.date_input("Availability Date")
    furnishing = st.selectbox("Furnishing Type", ["full", "semi", "unfurnished"])
    flooring_type = st.selectbox("Flooring Type", ["Marble", "Wood", "Tiles", "Vinyl", "Laminate"])
    roof_type = st.selectbox("Roof Type", ["Concrete", "Tile", "Metal", "Slate"])
    property_style = st.selectbox("Property Style", ["Modern", "Traditional", "Victorian", "Contemporary"])
    sewer_system = st.selectbox("Sewer System", ["Public", "Septic"])
with st.expander("Amenities"):
    security = st.multiselect("Security Features", ["Police Station Nearby", "CCTV", "Gated Community", "Security Guards"], default=["CCTV", "Gated Community"])
    gym = st.checkbox("Gym Available", value=True)
    swimming_pool = st.checkbox("Swimming Pool Available", value=False)
    convention_hall = st.checkbox("Convention Hall Available", value=False)
    parking = st.selectbox("Parking Availability", ["car", "bike", "both", "no parking"], index=2)
    property_conditions = st.selectbox("Current Property Condition", ["vacant", "tenant on notice", "new property"])
    bathrooms = st.number_input("Number of Bathrooms", min_value=1, step=1)
    balcony = st.number_input("Number of Balconies", min_value=0, step=1)
    water_supply = st.selectbox("Water Supply", ["corporation", "borewell", "both"])
    view = st.selectbox("View", ["Sea View", "City View", "Garden View", "Lake View", "None"])

```

```

lot_shape = st.selectbox("Lot Shape", ["Rectangular", "Irregular", "Square"])
with st.expander("Nearby property"):
    nearby_schools = st.text_input("Nearby Schools")
    nearby_hospitals = st.text_input("Nearby Hospitals")
    nearby_supermarkets = st.text_input("Nearby Supermarkets")
    it_hub_present = st.checkbox("IT Hub Present Nearby")
with st.expander("Photos of Property"):
    photos = st.text_input("Past the Property image links")
with st.expander("Financial Details"):
    deposit_cost = st.number_input("Deposit Cost", min_value=1000, step=1000, value=50000)
    hoa_fees = st.number_input("HOA Fees ", min_value=0, step=100, value=500)
    maint_fee = st.number_input("Maintenance Fee per Sqft", min_value=1, step=1, value=2)
with st.expander("Owner Details"):
    owner_name = st.text_input("Owner Name")
    owner_email = st.text_input("Owner Email")
    address = st.text_area("Full Address")
    pincode = st.text_input("Pincode")
    state = st.text_input("State", value="Tamil Nadu")
if st.button("Submit Property Data"):
    new_entry = {
        "City": city,
        "Locality": locality,
        "Property Purpose": property_purpose,
        "House Type": house_type,
        "BHK": bhk,
        "House Floor": house_floor,
        "Building Floor": building_floor,
        "Property Age": property_age,
        "Facing": facing,
        "Square Feet": sqft,
        "Total sqft Cost": sqft * sqft_cost,
        "Security": ", ".join(security),
        "Maintenance Fee": maint_fee * sqft,
        "Gym": gym,
        "Convention Hall": convention_hall,
        "Property Condition": property_conditions,
        "Bathrooms": bathrooms,
        "Balcony": balcony,
        "Water Supply": water_supply,
        "Parking": parking,
        "Availability Date": availability_date,
        "Deposit Cost": deposit_cost,
        "Furnishing": furnishing,
        "HOA Fees ": hoa_fees,
        "Flooring Type": flooring_type,
        "Roof Type": roof_type,
        "Property Style": property_style,
        "View": view,
        "Nearby Schools": nearby_schools,
    }

```

```

    "Nearby Hospitals": nearby_hospitals,
    "Nearby Supermarkets": nearby_supermarkets,
    "IT Hub Present": it_hub_present,
    "Address": address,
    "Owner_Name": owner_name,
    "Owner Email": owner_email,
    "Pincode": pincode,
    "State": state,
    "Swimming Pool": swimming_pool,
    "Photos": photos,
    "Lot Shape": lot_shape,
    "Sewer System": sewer_system
}

existing_records = df.drop(columns=["Property_ID"], errors="ignore")
if any(existing_records.apply(lambda row: row.to_dict() == new_entry, axis=1)):
    st.error("Property already exists in the database!")
else:
    property_id = generate_unique_property_id(df)
    new_entry["Property_ID"] = property_id
    df = pd.concat([df, pd.DataFrame([new_entry])], ignore_index=True)
    save_data_to_s3(df)
    st.success(f"Property added successfully! Generated Property ID: {property_id}")

```

Customer Code:

```

import streamlit as st
import boto3
import pandas as pd
import os
import openai
from openai import OpenAI
from datetime import datetime
from google import genai
# Instantiate the client
#                                     client = 
OpenAI(api_key="sk-proj-Jwz-F8GuN8o3bybXfkasJydTb6-esS9eQ0YuF--ULSLjOgfi9aFQBPCwdtTJR58I3K2We
yরpWZT3BlbkFJIBibmlBc4PcsE_7HuaCX7G0wChCy9sPB7NB2hMgTwAm1W8IXt7z8i_fHwdDZr4xDqMtJAip2
wA")
# Only run this block for Gemini Developer API
client = genai.Client(api_key='AIzaSyBY4Rj0QyMtK641PZBrwFTaqwRInHmqzEk')
# AWS credentials (Replace with st.secrets in production)
AWS_ACCESS_KEY = "AKIAYKFQRCZXOXDH4G6"
AWS_SECRET_KEY = "pANF8lHi/wtKg/D0UCJ3Ogc5juPhHirOK1tVzLuE"
AWS_REGION = "us-east-1"
# Initialize S3 client
s3_client = boto3.client(
    "s3",
    aws_access_key_id=AWS_ACCESS_KEY,

```

```

aws_secret_access_key=AWS_SECRET_KEY,
region_name=AWS_REGION,
)
# Reusable UI component
def get_common_inputs(is_rent=False):
    property_purpose = st.selectbox("Property Purpose", ["rent", "buy"])
    city = st.selectbox("City", ["chennai"])
    localities = ["Adyar", "Anna Nagar", "T. Nagar", "Velachery", "OMR", "Tambaram", "Porur", "Guindy",
                  "Mylapore", "Kilpauk", "Perambur", "Vadapalani", "Kodambakkam", "Royapettah", "Egmore",
                  "Pallavaram", "Avadi", "Chromepet", "Sholinganallur", "Ambattur"]
    locality = st.selectbox("Locality", localities)

    col1, col2 = st.columns(2)
    with col1:
        house_type = st.selectbox("House Type", ["apartment", "individual house", "villa", "gated community"])
        bkh = st.number_input("BHK", 1, 10, 2)
        house_floor = st.number_input("House Floor", 0, 50, 1)
        building_floor = st.number_input("Building Floor", 1, 50, 5)
        property_age = st.number_input("Property Age", 0, 100, 5)
        parking = st.selectbox("Parking", ["car", "bike", "both", "no parking"])
    with col2:
        facing = st.selectbox("Facing", ["north", "south", "east", "west", "north-east", "north-west", "south-east",
                                         "south-west"])
        sqft = st.number_input("Square Feet", 100, 10000, 1000)
        security = st.selectbox("Security", ["none", "cctv", "guards", "gated"])
        gym = st.selectbox("Gym", ["yes", "no"])
        convention = st.selectbox("Convention Hall", ["yes", "no"])
        water = st.selectbox("Water Supply", ["corporation", "borewell", "both"])
        bathrooms = st.slider("Bathrooms", 1, 5, 2)
        balcony = st.slider("Balcony", 0, 5, 1)
        schools = st.slider("Nearby Schools", 0, 5, 2)
    col3, col4 = st.columns(2)
    with col3:
        availability = st.date_input("Availability Date")
        furnishing = st.selectbox("Furnishing", ["full", "semi", "unfurnished"])
        flooring = st.selectbox("Flooring Type", ["marble", "wood", "tiles", "vinyl", "laminate"])
        roof = st.selectbox("Roof Type", ["concrete", "tile", "metal", "slate"])
        style = st.selectbox("Property Style", ["modern", "traditional", "victorian", "contemporary"])
    with col4:
        it_hub = st.selectbox("IT Hub Present", ["none", "tcs", "infosys", "cts", "zoho", "accenture", "wipro"])
        pool = st.selectbox("Swimming Pool", ["yes", "no"])
        hospitals = st.slider("Nearby Hospitals", 0, 5, 1)
        supermarkets = st.slider("Nearby Supermarkets", 0, 5, 2)
    inputs = {
        "property_purpose": property_purpose.lower(),
        "locality": locality.lower(),
        "house type": house_type.lower(),
        "bkh": bkh,
        "house floor": house_floor,

```

```

    "building floor": building_floor,
    "property age": property_age,
    "facing": facing.lower(),
    "square feet": sqft,
    "security": security.lower(),
    "gym": gym.lower(),
    "convention hall": convention.lower(),
    "parking": parking.lower(),
    "water supply": water.lower(),
    "bathrooms": bathrooms,
    "balcony": balcony,
    "availability date": str(availability),
    "furnishing": furnishing.lower(),
    "flooring type": flooring.lower(),
    "roof type": roof.lower(),
    "property style": style.lower(),
    "nearby schools": schools,
    "nearby hospitals": hospitals,
    "nearby supermarkets": supermarkets,
    "it hub present": it_hub.lower(),
    "swimming pool": pool.lower()
}
return inputs
def fetch_properties(bucket_name, file_key, filters):
    try:
        response = s3_client.get_object(Bucket=bucket_name, Key=file_key)
        df = pd.read_csv(response["Body"])
        df.columns = [col.strip().lower().replace(" ", "_") for col in df.columns]
        # Simplified filter keys
        normalized_filters = {
            k.strip().lower().replace(" ", "_"): v for k, v in filters.items()
        }
        # Manually filter with debug logging
        debug_keys = ["property_purpose", "locality", "house_type", "bhk", "property_age"]
        for key in debug_keys:
            if key in df.columns and key in normalized_filters:
                value = normalized_filters[key]
                st.write(f" Filtering: {key} == {value}")
                if isinstance(value, str):
                    df = df[df[key].str.strip().str.lower() == value.strip().lower()]
                else:
                    df = df[df[key] == value]
                st.write(" Filtered Preview", df.head())
        return df
    except Exception as e:
        st.error(f"Error fetching file from S3: {e}")
        return pd.DataFrame()
def get_gemini_insight(df):
    if df.empty:

```

```

    return "No data to analyze."
summary_data = df.head(10).to_string(index=False)
prompt = f"""
You are a real estate expert.
Analyze the following property listings and recommend the best one, including clear reasoning and give them a
bargain amount or the percentage.

==== Property Data ====
{summary_data}
=====

Please follow this format for the answer:
Property Recommendation:
[Your recommendation here]
Reasoning:
1. [Point 1]
2. [Point 2]
3. [Point 3]
...
Concerns:
1. [Concern 1]
2. [Concern 2]
...
Negotiation Advice:
1. [Advice 1]
2. [Advice 2]
3. [Advice 3]
...
Important Considerations Before Finalizing:
1. [Consideration 1]
2. [Consideration 2]
3. [Consideration 3]

Write the response in plain text.
Avoid using stars (*), dashes (-), or any bullet symbols.
Make the text look like a clean written report.

"""

try:
    response = client.models.generate_content(
        model='gemini-2.0-flash-001',
        contents=prompt
    )
    return response.text.strip()
except Exception as e:
    return f"Gemini Error: {e}"
# ┌ Main App
def main():
    st.title("Real Estate Property Insights for Customers")
    bucket_name = "magedataset"
    file_key = "real_state_data.csv"
    filters = get_common_inputs(is_rent=True)
    if st.button("Search Properties"):

```

```

result_df = fetch_properties(bucket_name, file_key, filters)
if not result_df.empty:
    st.success(f"⚡️Found {len(result_df)} matching properties.")
    st.dataframe(result_df)
    insight = get_gemini_insight(result_df)
    st.text_area("🔗 OpenAI Recommendation", insight, height=200)
else:
    st.warning("⚠️No matching properties found.")
if __name__ == "__main__":
    main()

```

Owner Code:

```

import streamlit as st
import pandas as pd
import openai
import os
from openai import OpenAI
from pyspark.sql import SparkSession
from pyspark.sql.functions import when, col, months_between, current_date, lower
from pyspark.ml import PipelineModel
from pyspark.ml.regression import LinearRegressionModel
from pyspark.sql import Row
from google import genai
# ----- Setup -----
# Spark config
os.environ["PYSPARK_PYTHON"] =
"/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/myenv/bin/python"
os.environ["PYSPARK_DRIVER_PYTHON"] =
"/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/myenv/bin/python"
spark = SparkSession.builder.appName("RealEstateApp").getOrCreate()
# Load models
buy_pipeline = PipelineModel.load("models/preprocessing_pipeline/content/preprocessing_pipeline")
buy_model = LinearRegressionModel.load("models/house_price_model/content/house_price_model")
rent_pipeline = PipelineModel.load("models/rental_pipeline/content/rental_pipeline")
rent_model = LinearRegressionModel.load("models/rental_model/content/rental_model")
# OpenAI Key
openai.api_key = "key"
# ----- Streamlit UI -----
st.set_page_config(page_title="PropIntel", page_icon="🔗", layout="wide")
st.title("🔗 PropIntel: Real Estate Prediction")
option = st.radio("What would you like to predict?", ["🔗 Buy House Price", "🔗 Rent Monthly Cost"]) #
----- OpenAI Insight -----
client = OpenAI(api_key=openai.api_key)
def get_openai_insight(user_input_dict, prediction_value, mode="buy"):
    input_summary = "\n".join([f'{k}: {v}' for k, v in user_input_dict.items()])
    prediction_line = (
        f"Predicted Buy Price per Sqft: ₹{prediction_value}, " if mode == "buy"
        else f"Estimated Monthly Rent: ₹{prediction_value}")

```

```

)
prompt = f"""
A owner provided the following real estate property details:
{input_summary}
{prediction_line}
you need to say "Great!!" and you need to appriciate
Based on this information, provide a helpful real estate suggestion for increase in price or not if possible provide
the numbers,
you need to make the
owner think its the best and updated on going sale market trend. Also give the owner profit of this choice..
"""

chat = client.chat.completions.create(
    model="gpt-3.5-turbo",
    messages=[{"role": "user", "content": prompt}],
    max_tokens=500,
    temperature=0.7
)
return chat.choices[0].message.content.strip()
# ----- Form -----
def get_common_inputs(is_rent=False):
    city = st.selectbox("City", ["chennai"])
    localities = ["Adyar", "Anna Nagar", "T. Nagar", "Velachery", "OMR", "Tambaram", "Porur", "Guindy",
                  "Mylapore", "Kilpauk", "Perambur", "Vadapalani", "Kodambakkam", "Royapettah", "Egmore",
                  "Pallavaram", "Avadi", "Chromepet", "Sholinganallur", "Ambattur"]
    locality = st.selectbox("Locality", localities)
    col1, col2 = st.columns(2)
    with col1:
        house_type = st.selectbox("House Type", ["apartment", "individual house", "villa", "gated community"])
        bhk = st.number_input("BHK", 1, 10, 2)
        house_floor = st.number_input("House Floor", 0, 50, 1)
        building_floor = st.number_input("Building Floor", 1, 50, 5)
        property_age = st.number_input("Property Age", 0, 100, 5)
        facing = st.selectbox("Facing", ["north", "south", "east", "west", "north-east", "north-west", "south-east",
                                         "south-west"])
        sqft = st.number_input("Square Feet", 100, 10000, 1000)
        bathrooms = st.slider("Bathrooms", 1, 5, 2)
        balcony = st.slider("Balcony", 0, 5, 1)
        lot = st.selectbox("Lot Shape", ["rectangular", "irregular", "square"])
        sewer = st.selectbox("Sewer System", ["public", "septic"])
    with col2:
        security = st.selectbox("Security", ["none", "cctv", "guards", "gated"])
        gym = st.selectbox("Gym", ["yes", "no"])
        pool = st.selectbox("Swimming Pool", ["yes", "no"])
        convention = st.selectbox("Convention Hall", ["yes", "no"])
        parking = st.selectbox("Parking", ["car", "bike", "both", "no parking"])
        water = st.selectbox("Water Supply", ["corporation", "borewell", "both"])
        flooring = st.selectbox("Flooring Type", ["marble", "wood", "tiles", "vinyl", "laminate"])
        roof = st.selectbox("Roof Type", ["concrete", "tile", "metal", "slate"])
        style = st.selectbox("Property Style", ["modern", "traditional", "victorian", "contemporary"])

```

```

view = st.selectbox("View", ["sea view", "city view", "garden view", "lake view", "none"])
availability = st.date_input("Availability Date")
furnishing = st.selectbox("Furnishing", ["full", "semi", "unfurnished"])
hoa = st.number_input("HOA Fees", 0.0, 100000.0, 1500.0)
schools = st.slider("Nearby Schools", 0, 5, 2)
hospitals = st.slider("Nearby Hospitals", 0, 5, 1)
supermarkets = st.slider("Nearby Supermarkets", 0, 5, 2)
it_hub = st.selectbox("IT Hub Present", ["none", "tcs", "infosys", "cts", "zoho", "accenture", "wipro"])
photos = st.selectbox("Photos Available?", ["yes", "no"])
inputs = {
    "City": city,
    "Locality": locality,
    "House Type": house_type,
    "BHK": bhk,
    "House Floor": house_floor,
    "Building Floor": building_floor,
    "Property Age": property_age,
    "Facing": facing.lower(),
    "Square Feet": sqft,
    "Security": security.lower(),
    "Gym": gym.lower(),
    "Convention Hall": convention.lower(),
    "Parking": parking.lower(),
    "Water Supply": water.lower(),
    "Bathrooms": bathrooms,
    "Balcony": balcony,
    "Availability Date": str(availability),
    "Furnishing": furnishing.lower(),
    "HOA Fees": hoa,
    "Flooring Type": flooring.lower(),
    "Roof Type": roof.lower(),
    "Property Style": style.lower(),
    "View": view.lower(),
    "Nearby Schools": schools,
    "Nearby Hospitals": hospitals,
    "Nearby Supermarkets": supermarkets,
    "IT Hub Present": it_hub.lower(),
    "Swimming Pool": pool.lower(),
    "Lot Shape": lot.lower(),
    "Sewer System": sewer.lower(),
    "Photos": photos.lower()
}
if is_rent:
    deposit = st.number_input("Deposit Cost", 0, 1000000, 50000)
    hoa_r = st.number_input("HOA Fees ", 0.0, 100000.0, 1500.0)
    inputs["Deposit Cost"] = deposit
    inputs["HOA Fees "] = hoa_r
return inputs
# ----- Buy Prediction -----

```

```

if option == "Buy House Price":
    st.header("Buy House Price Prediction")
    user_input = get_common_inputs()
    if st.button("Predict Buy Price"):
        user_row = Row(**{**user_input, "Property Purpose": "buy"})
        df = spark.createDataFrame([user_row])
        df = df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
        df = df.withColumn("Months_To_Availability", months_between(col("Availability Date"), current_date()))
        df = df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
        df = df.withColumn("Amenity_Score", (
            when(lower(col("Gym")) == "yes", 1).otherwise(0) +
            when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
            when(lower(col("Security")) != "none", 1).otherwise(0) +
            when(lower(col("Convention Hall")) == "yes", 1).otherwise(0)
        ))
        df = df.withColumn("Is_Near_IT_Hub", when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
        transformed = buy_pipeline.transform(df)
        prediction = buy_model.transform(transformed)
        prediction = prediction.withColumn("Total_Predicted_Price", col("prediction") * col("Square Feet"))
        row = prediction.select("prediction", "Total_Predicted_Price").first()
        price_per_sqft = round(row["prediction"])
        total_price = round(row["Total_Predicted_Price"])
        st.success(f"Predicted Price per Sqft: ₹{price_per_sqft},")
        st.success(f"Total Estimated Property Price: ₹{total_price},")
        insight = get_openai_insight(user_input, price_per_sqft, "buy")
        st.text_area("AI Insight", insight, height=300)
# ----- Rent Prediction -----
elif option == "Rent Monthly Cost":
    st.header("Rent Monthly Cost Prediction")
    user_input = get_common_inputs(is_rent=True)
    if st.button("Predict Rent Price"):
        user_row = Row(**{**user_input, "Property Purpose": "rent"})
        df = spark.createDataFrame([user_row])
        df = df.withColumn("BHK_Sqft", col("BHK") * col("Square Feet"))
        df = df.withColumn("Months_To_Availability", months_between(col("Availability Date"), current_date()))
        df = df.withColumn("Has_Photos", when(col("Photos").isNotNull(), 1).otherwise(0))
        df = df.withColumn("Amenity_Score", (
            when(lower(col("Gym")) == "yes", 1).otherwise(0) +
            when(lower(col("Swimming Pool")) == "yes", 1).otherwise(0) +
            when(lower(col("Security")) != "none", 1).otherwise(0) +
            when(lower(col("Convention Hall")) == "yes", 1).otherwise(0)
        ))
        df = df.withColumn("Is_Near_IT_Hub", when(lower(col("IT Hub Present")) != "none", 1).otherwise(0))
        transformed = rent_pipeline.transform(df)
        prediction = rent_model.transform(transformed)
        rent = round(prediction.select("prediction").first()["prediction"])
        st.success(f"Estimated Monthly Rent: ₹{rent},")
        insight = get_openai_insight(user_input, rent, "rent")
        st.text_area("AI Insight", insight, height=300)

```

Log:

2025-03-14 21:48:42 [INFO] Question: how many number of row are there in the dataset?
2025-03-14 21:48:42 [INFO] Running PandasAI with openai LLM...
2025-03-14 21:48:42 [INFO] Prompt ID: 4956f985-8505-40fe-bff5-630f8afc51c9
2025-03-14 21:48:42 [INFO] Executing Pipeline: GenerateChatPipeline
2025-03-14 21:48:43 [INFO] Executing Step 0: ValidatePipelineInput
2025-03-14 21:48:43 [INFO] Executing Step 1: CacheLookup
2025-03-14 21:48:43 [INFO] Executing Step 2: PromptGeneration
2025-03-14 21:48:43 [INFO] Using prompt: <dataframe>
dfs[0]:2999999x5
1,Melissa Santos,metacin,male,503
1377984,Kristopher Chan,avil,female,794
1457725,James Cisneros,paracetamol,male,548
857167,Clayton Saunders,metacin,female,986
</dataframe>
Update this initial code:
```python  
# TODO: import the required dependencies  
import pandas as pd  
# Write code here  
# Declare result var:  
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest\_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp\_chart.png" }  
```  
QUERY
how many number of row are there in the dataset?
Variable `dfs: list[pd.DataFrame]` is already declared.
At the end, declare "result" variable as a dictionary of type and value.
If you are asked to plot a chart, use "matplotlib" for charts, save as png.
Generate python code and return full updated code:
2025-03-14 21:48:43 [INFO] Executing Step 3: CodeGenerator
2025-03-14 21:48:43 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 401 Unauthorized"
2025-03-14 21:48:43 [ERROR] Pipeline failed on step 3: **Error** code: 401 - {'error': {'message': 'Incorrect API key provided: API KEY. You can find your API key at https://platform.openai.com/account/api-keys.', 'type': 'invalid_request_error', 'param': None, 'code': 'invalid_api_key'}}
2025-03-14 21:50:09 [INFO] Question: total number of sales?
2025-03-14 21:50:09 [INFO] Running PandasAI with openai LLM...
2025-03-14 21:50:09 [INFO] Prompt ID: 15bc3d1a-23ea-40be-8766-948917730845
2025-03-14 21:50:09 [INFO] Executing Pipeline: GenerateChatPipeline
2025-03-14 21:50:09 [INFO] Executing Step 0: ValidatePipelineInput
2025-03-14 21:50:09 [INFO] Executing Step 1: CacheLookup
2025-03-14 21:50:09 [INFO] Executing Step 2: PromptGeneration
2025-03-14 21:50:09 [INFO] Using prompt: <dataframe>
dfs[0]:51290x16
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method

2018-11-07,20:48:51,2.0,98188,Male,Mobile,First
 SignUp,Home &
 Furniture,Mixer/Juicer,196.0,2.0,,55.6.,Medium,not_defined
 2018-06-17,19:47:39,,70794,Male,Web,Guest,Electronic,Fans,224.0,,0.2,27.1,3.8.,e_wallet
 2018-11-18,17:46:06,10.5,72004,Female,Web,New ,Auto &
 Accessories,Jeans,,5.0,0.1,51.7,5.4,Critical,money_order
 </dataframe>
 Update this initial code:
 ```python  
 # TODO: import the required dependencies  
 import pandas as pd  
 # Write code here  
 # Declare result var:  
 type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest\_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp\_chart.png" }  
 ```  
QUERY
 total number of sales?
 Variable `dfs: list[pd.DataFrame]` is already declared.
 At the end, declare "result" variable as a dictionary of type and value.
 If you are asked to plot a chart, use "matplotlib" for charts, save as png.
 Generate python code and return full updated code:
 2025-03-14 21:50:09 [INFO] Executing Step 3: CodeGenerator
 2025-03-14 21:50:09 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 401 Unauthorized"
 2025-03-14 21:50:09 [ERROR] Pipeline failed on step 3: **Error** code: 401 - {'error': {'message': 'Incorrect API key provided: sk-proj-*****FKWY. You can find your API key at https://platform.openai.com/account/api-keys.'}, 'type': 'invalid_request_error', 'param': None, 'code': 'invalid_api_key'}
 2025-03-14 21:51:57 [INFO] Question: total number of sales?
 2025-03-14 21:51:57 [INFO] Running PandasAI with openai LLM...
 2025-03-14 21:51:57 [INFO] Prompt ID: e9474b23-bfe5-4fd7-91bd-2d9ab23e8f15
 2025-03-14 21:51:57 [INFO] Executing Pipeline: GenerateChatPipeline
 2025-03-14 21:51:57 [INFO] Executing Step 0: ValidatePipelineInput
 2025-03-14 21:51:57 [INFO] Executing Step 1: CacheLookup
 2025-03-14 21:51:57 [INFO] Executing Step 2: PromptGeneration
 2025-03-14 21:51:57 [INFO] Using prompt: <dataframe>
 dfs[0]:51290x16
 Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
 Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
 2018-12-28,20:22:17,9.0,31244,Male,Mobile,Guest,Fashion,Jeans,,0.4,123.3,,Low,not_defined
 2018-11-12,13:40:52,2.0,72810,Female,Web,New ,Auto & Accessories,Fossil Watch,65.0,4.0,,0.5,7.3,,money_order
 2018-01-16,17:45:41,,20141,Female,Mobile,Member,Home &
 Furniture,Shirts,248.0,5.0,0.3,107.0,9.4,Medium,e_wallet
 </dataframe>
 Update this initial code:
 ```python  
 # TODO: import the required dependencies  
 import pandas as pd

```

Write code here
Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```
### QUERY
total number of sales?
Variable `dfs: list[pd.DataFrame]` is already declared.
At the end, declare "result" variable as a dictionary of type and value.
If you are asked to plot a chart, use "matplotlib" for charts, save as png.
Generate python code and return full updated code:
2025-03-14 21:51:57 [INFO] Executing Step 3: CodeGenerator
2025-03-14 21:51:58 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 429 Too Many Requests"
2025-03-14 21:51:58 [INFO] Retrying request to /chat/completions in 0.444971 seconds
2025-03-14 21:51:59 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 429 Too Many Requests"
2025-03-14 21:51:59 [INFO] Retrying request to /chat/completions in 0.839569 seconds
2025-03-14 21:52:00 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 429 Too Many Requests"
2025-03-14 21:52:00 [ERROR] Pipeline failed on step 3: Error code: 429 - {'error': {'message': 'You exceeded your current quota, please check your plan and billing details. For more information on this error, read the docs: https://platform.openai.com/docs/guides/error-codes/api-errors.'}, 'type': 'insufficient_quota', 'param': None, 'code': 'insufficient_quota'}
2025-03-14 21:56:54 [INFO] Question: total number of sales?
2025-03-14 21:56:54 [INFO] Running PandasAI with openai LLM...
2025-03-14 21:56:54 [INFO] Prompt ID: 248e0a5d-5eaf-4f2a-8f4f-c935ceb85499
2025-03-14 21:56:54 [INFO] Executing Pipeline: GenerateChatPipeline
2025-03-14 21:56:54 [INFO] Executing Step 0: ValidatePipelineInput
2025-03-14 21:56:54 [INFO] Executing Step 1: CacheLookup
2025-03-14 21:56:54 [INFO] Executing Step 2: PromptGeneration
2025-03-14 21:56:54 [INFO] Using prompt: <dataframe>
dfs[0]:51290x16
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
2018-02-16,14:30:31,,96925,Female,Mobile,First SignUp,Home & Furniture,Curtains,,3.0,0.4,116.6,5.3,,credit_card
2018-08-24,18:49:35,8.0,0.91460,Male,Mobile,Member,Fashion,Fossil Watch,85.0,2.0,,39.0,9.7,Critical,money_order
2018-01-07,12:31:21,1.0,0.63372,Female,Web,Guest,Electronic,Dinning Tables,104.0,,0.3,35.9,,Medium,not_defined
</dataframe>
Update this initial code:
```
python
TODO: import the required dependencies
import pandas as pd
Write code here
Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```

```

```

### ### QUERY

total number of sales?

Variable `dfs: list[pd.DataFrame]` is already declared.

At the end, declare "result" variable as a dictionary of type and value.

If you are asked to plot a chart, use "matplotlib" for charts, save as png.

Generate python code and return full updated code:

2025-03-14 21:56:54 [INFO] Executing Step 3: CodeGenerator

2025-03-14 21:56:55 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 401 Unauthorized"

2025-03-14 21:56:55 [ERROR] Pipeline failed on step 3: **Error** code: 401 - {'error': {'message': 'Incorrect API key provided'}}

sk-proj-\*\*\*\*\*  
\*\*\*\*\*NIFU. You can find your API key at  
<https://platform.openai.com/account/api-keys>.', 'type': 'invalid\_request\_error', 'param': None, 'code': 'invalid\_api\_key'}}

2025-03-14 21:57:58 [INFO] Question: total sales?

2025-03-14 21:57:58 [INFO] Running PandasAI with openai LLM...

2025-03-14 21:57:58 [INFO] Prompt ID: 516b91d6-5870-45b3-bd8b-68ae729ab497

2025-03-14 21:57:58 [INFO] Executing Pipeline: GenerateChatPipeline

2025-03-14 21:57:58 [INFO] Executing Step 0: ValidatePipelineInput

2025-03-14 21:57:58 [INFO] Executing Step 1: CacheLookup

2025-03-14 21:57:58 [INFO] Executing Step 2: PromptGeneration

2025-03-14 21:57:58 [INFO] Using prompt: <dataframe>

dfs[0]:51290x16

Order\_Date,Time,Aging,Customer\_Id,Gender,Device\_Type,Customer\_Login\_type,Product\_Category,Product,Sales,Quantity,Discount,Profit,Shipping\_Cost,Order\_Priority,Payment\_method

2018-06-09,11:39:02,1.0,71421,Female,Mobile,First SignUp,Electronic,Apple Laptop,,,110.4,8.9,,debit\_card

2018-01-21,12:07:26,2.0,19201,Male,Mobile,New ,Auto & Accessories,Titak watch,65.0,3.0,0.1,9.3,,Low,money\_order

2018-05-09,13:26:08,,86399,Male,Web,Member,Home

&

Furniture,Sofas,220.0,4.0,0.2,131.7,1.7,Critical,not\_defined

</dataframe>

Update this initial code:

```python

TODO: import the required dependencies

import pandas as pd

Write code here

Declare result var:

type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }

```

### ### QUERY

total sales?

Variable `dfs: list[pd.DataFrame]` is already declared.

At the end, declare "result" variable as a dictionary of type and value.

If you are asked to plot a chart, use "matplotlib" for charts, save as png.

Generate python code and return full updated code:

2025-03-14 21:57:58 [INFO] Executing Step 3: CodeGenerator

2025-03-14 21:57:58 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 401 Unauthorized"

2025-03-14 21:57:58 [ERROR] Pipeline failed on step 3: **Error** code: 401 - {'error': {'message': 'Incorrect API key provided:  
sk-proj\_\*\*\*\*\*  
\*\*\*\*\*NIFU. You can find your API key at  
<https://platform.openai.com/account/api-keys>.', 'type': 'invalid\_request\_error', 'param': None, 'code': 'invalid\_api\_key'}}

2025-03-14 22:10:41 [INFO] Question: total sales?

2025-03-14 22:10:41 [INFO] Running PandasAI with openai LLM...

2025-03-14 22:10:41 [INFO] Prompt ID: 4063a424-bffd-4dfb-8817-0e84c6c9242f

2025-03-14 22:10:41 [INFO] Executing Pipeline: GenerateChatPipeline

2025-03-14 22:10:41 [INFO] Executing Step 0: ValidatePipelineInput

2025-03-14 22:10:41 [INFO] Executing Step 1: CacheLookup

2025-03-14 22:10:41 [INFO] Executing Step 2: PromptGeneration

2025-03-14 22:10:41 [INFO] Using prompt: <dataframe>

dfs[0]:51290x16

Order\_Date,Time,Aging,Customer\_Id,Gender,Device\_Type,Customer\_Login\_type,Product\_Category,Product,Sales,  
Quantity,Discount,Profit,Shipping\_Cost,Order\_Priority,Payment\_method

2018-01-08,17:01:39,2.0,0.95235,Female,Web,New ,Home & Furniture,Formal Shoes,,3.0,0.4,64.7,,Low,credit\_card

2018-09-28,23:36:17,1.0,0.97846,Female,Mobile,Guest,Electronic,Keyboard,54.0,5.0,,25.3,2.9,Medium,not\_defined

2018-10-25,13:11:29,,33392,Male,Mobile,Member,Fashion,Car                                  Pillow                          &                          Neck  
Rest,104.0,,0.2,137.6,8.8,,money\_order

</dataframe>

Update this initial code:

```
```python
# TODO: import the required dependencies
import pandas as pd
# Write code here
# Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```

```

### QUERY

total sales?

Variable `dfs: list[pd.DataFrame]` is already declared.

At the end, declare "result" variable as a dictionary of type and value.

If you are asked to plot a chart, use "matplotlib" for charts, save as png.

Generate python code and return full updated code:

2025-03-14 22:10:41 [INFO] Executing Step 3: CodeGenerator

2025-03-14 22:10:44 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 200 OK"

2025-03-14 22:10:44 [INFO] Prompt used:

<dataframe>

dfs[0]:51290x16

Order\_Date,Time,Aging,Customer\_Id,Gender,Device\_Type,Customer\_Login\_type,Product\_Category,Product,Sales,  
Quantity,Discount,Profit,Shipping\_Cost,Order\_Priority,Payment\_method

```
2018-01-08,17:01:39,2.0,95235,Female,Web,New ,Home & Furniture,Formal Shoes,,3.0,0.4,64.7,,Low,credit_card
2018-09-28,23:36:17,1.0,97846,Female,Mobile,Guest,Electronic,Keyboard,54.0,5.0,,25.3,2.9,Medium,not_defined
2018-10-25,13:11:29,,33392,Male,Mobile,Member,Fashion,Car Pillow & Neck
Rest,104.0.,0.2,137.6,8.8.,money_order
```

```
</dataframe>
```

Update this initial code:

```
'''python
TODO: import the required dependencies
import pandas as pd
Write code here
Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest
salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value":
pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
'''
```

### QUERY

total sales?

Variable `dfs: list[pd.DataFrame]` is already declared.

At the end, declare "result" variable as a dictionary of type and value.

If you are asked to plot a chart, use "matplotlib" for charts, save as png.

Generate python code and return full updated code:

```
2025-03-14 22:10:44 [INFO] Code generated:
```

```
'''
```

```
TODO: import the required dependencies
import pandas as pd
Assuming dfs is already defined and contains the DataFrames
For this example, we will use the first DataFrame in the list
df = dfs[0]
Convert the 'Sales' column to numeric, handling any potential errors
df['Sales'] = pd.to_numeric(df['Sales'], errors='coerce')
Calculate total sales
total_sales = df['Sales'].sum()
Declare result var:
result = {
 "type": "number",
 "value": total_sales
}
'''
```

```
2025-03-14 22:10:44 [INFO] Executing Step 4: CachePopulation
```

```
2025-03-14 22:10:44 [INFO] Executing Step 5: CodeCleaning
```

```
2025-03-14 22:10:44 [INFO]
```

Code running:

```
'''
```

```
df = dfs[0]
df['Sales'] = pd.to_numeric(df['Sales'], errors='coerce')
total_sales = df['Sales'].sum()
result = {'type': 'number', 'value': total_sales}
```

```

2025-03-14 22:10:44 [INFO] Executing Step 6: CodeExecution

2025-03-14 22:10:44 [ERROR] Failed with **error**: Traceback (most recent call last):

```
File "/opt/anaconda3/lib/python3.12/site-packages/pandasai/pipelines/chat/code_execution.py", line 85, in execute
  result = self.execute_code(code_to_run, code_context)
           ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
```

```
  File  "/opt/anaconda3/lib/python3.12/site-packages/pandasai/pipelines/chat/code_execution.py", line 174, in
execute_code
```

```
    exec(code, environment)
```

```
  File "<string>", line 2, in <module>
```

```
    File  "/opt/anaconda3/lib/python3.12/site-packages/pandasai/safe_libs/restricted_pandas.py", line 109, in
__getattr__
```

```
      raise AttributeError(f"'{name}' is not allowed in RestrictedPandas")
```

```
AttributeError: 'to_numeric' is not allowed in RestrictedPandas
```

2025-03-14 22:10:44 [WARNING] Failed to execute code retrying with a correction framework [retry number: 1]

2025-03-14 22:10:44 [INFO] Executing Pipeline: ErrorCorrectionPipeline

2025-03-14 22:10:44 [INFO] Executing Step 0: ErrorPromptGeneration

2025-03-14 22:10:44 [INFO] Using prompt: <dataframe>

```
dfs[0]:51290x16
```

```
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
```

```
2018-01-08,17:01:39,2.0,0.95235,Female,Web,New ,Home & Furniture,Formal Shoes,,3.0,0.4,64.7,,Low,credit_card
```

```
2018-09-28,23:36:17,1.0,0.97846,Female,Mobile,Guest,Electronic,Keyboard,54.0,5.0,,25.3,2.9,Medium,not_defined
```

```
2018-10-25,13:11:29,,33392,Male,Mobile,Member,Fashion,Car          Pillow          &          Neck
```

```
Rest,104.0,,0.2,137.6,8.8,,money_order
```

```
</dataframe>
```

The user asked the following question:

QUERY

total sales?

You generated this python code:

```
df = dfs[0]
df['Sales'] = pd.to_numeric(df['Sales'], errors='coerce')
total_sales = df['Sales'].sum()
result = {'type': 'number', 'value': total_sales}
```

It fails with the following **error**:

Traceback (most recent call last):

```
File "/opt/anaconda3/lib/python3.12/site-packages/pandasai/pipelines/chat/code_execution.py", line 85, in execute
  result = self.execute_code(code_to_run, code_context)
           ^^^^^^^^^^^^^^^^^^^^^^^^^^
```

```
  File  "/opt/anaconda3/lib/python3.12/site-packages/pandasai/pipelines/chat/code_execution.py", line 174, in
execute_code
```

```
    exec(code, environment)
```

```
  File "<string>", line 2, in <module>
```

```
    File  "/opt/anaconda3/lib/python3.12/site-packages/pandasai/safe_libs/restricted_pandas.py", line 109, in
__getattr__
```

```
      raise AttributeError(f"'{name}' is not allowed in RestrictedPandas")
```

```
AttributeError: 'to_numeric' is not allowed in RestrictedPandas
```

Fix the python code above and return the new python code:

2025-03-14 22:10:44 [INFO] Executing Step 1: CodeGenerator

2025-03-14 22:10:47 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 200 OK"

2025-03-14 22:10:47 [INFO] Prompt used:

<dataframe>

dfs[0]:51290x16

Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,

Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method

2018-01-08,17:01:39,2.0,0.95235,Female,Web,New ,Home & Furniture,Formal Shoes,,3.0,0.4,64.7,,Low,credit_card

2018-09-28,23:36:17,1.0,0.97846,Female,Mobile,Guest,Electronic,Keyboard,54.0,5.0,,25.3,2.9,Medium,not_defined

2018-10-25,13:11:29,,33392,Male,Mobile,Member,Fashion,Car Pillow & Neck

Rest,104.0,,0.2,137.6,8.8,,money_order

</dataframe>

The user asked the following question:

QUERY

total sales?

You generated this python code:

```
df = dfs[0]
```

```
df['Sales'] = pd.to_numeric(df['Sales'], errors='coerce')
```

```
total_sales = df['Sales'].sum()
```

```
result = {'type': 'number', 'value': total_sales}
```

It fails with the following **error**:

Traceback (most recent call last):

```
File "/opt/anaconda3/lib/python3.12/site-packages/pandasai/pipelines/chat/code_execution.py", line 85, in execute
    result = self.execute_code(code_to_run, code_context)
           ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
```

```
File "/opt/anaconda3/lib/python3.12/site-packages/pandasai/pipelines/chat/code_execution.py", line 174, in
execute_code
```

```
    exec(code, environment)
```

```
File "<string>", line 2, in <module>
```

```
    File "/opt/anaconda3/lib/python3.12/site-packages/pandasai/safe_libs/restricted_pandas.py", line 109, in
__getattr__
```

```
        raise AttributeError(f"'{name}' is not allowed in RestrictedPandas")
```

```
AttributeError: 'to_numeric' is not allowed in RestrictedPandas
```

Fix the python code above and return the new python code:

2025-03-14 22:10:47 [INFO] Code generated:

```

```
df = dfs[0]
```

```
Replace NaN values in 'Sales' with 0 for summation
```

```
df['Sales'] = df['Sales'].fillna(0)
```

```
Convert 'Sales' to float type
```

```
df['Sales'] = df['Sales'].astype(float)
```

```
total_sales = df['Sales'].sum()
```

```
result = {'type': 'number', 'value': total_sales}
```

```

2025-03-14 22:10:47 [INFO] Executing Step 2: CodeCleaning

2025-03-14 22:10:47 [INFO]

Code running:

```

```
df = dfs[0]
df['Sales'] = df['Sales'].fillna(0)
df['Sales'] = df['Sales'].astype(float)
total_sales = df['Sales'].sum()
result = {'type': 'number', 'value': total_sales}
```

2025-03-14 22:10:47 [INFO] Executing Step 7: ResultValidation
2025-03-14 22:10:47 [INFO] Answer: {'type': 'number', 'value': 7813411.0}
2025-03-14 22:10:47 [INFO] Executing Step 8: ResultParsing
2025-03-14 22:11:50 [INFO] Question: create a pie chart on customer login type column
2025-03-14 22:11:50 [INFO] Running PandasAI with openai LLM...
2025-03-14 22:11:50 [INFO] Prompt ID: 273b3ae5-5256-4e92-a88b-69d33233f881
2025-03-14 22:11:50 [INFO] Executing Pipeline: GenerateChatPipeline
2025-03-14 22:11:50 [INFO] Executing Step 0: ValidatePipelineInput
2025-03-14 22:11:50 [INFO] Executing Step 1: CacheLookup
2025-03-14 22:11:50 [INFO] Executing Step 2: PromptGeneration
2025-03-14 22:11:50 [INFO] Using prompt: <dataframe>
dfs[0]:51290x16
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
2018-09-29,22:13:43,,13996,Female,Web,First           SignUp,Home      &      Furniture,Fossil
Watch,62.0,1.0,0.4,39.6,12.6,High,not_defined
2018-01-06,21:06:53,9.0,55476,Male,Mobile,New ,Auto & Accessories,Iron,,4.0,0.1,34.4,,e_wallet
2018-02-21,14:51:48,8.0,44850,Male,Web,Member,Electronic,Car           &      Bike
Care,192.0,,,162.5,1.9,Medium,money_order
</dataframe>
Update this initial code:
```python
TODO: import the required dependencies
import pandas as pd
Write code here
Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```
### QUERY
create a pie chart on customer login type column
Variable `dfs: list[pd.DataFrame]` is already declared.
At the end, declare "result" variable as a dictionary of type and value.
If you are asked to plot a chart, use "matplotlib" for charts, save as png.
Generate python code and return full updated code:
2025-03-14 22:11:50 [INFO] Executing Step 3: CodeGenerator
2025-03-14 22:11:55 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 200 OK"
2025-03-14 22:11:55 [INFO] Prompt used:
<dataframe>
dfs[0]:51290x16

```

```

Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
2018-09-29,22:13:43,,13996,Female,Web,First           SignUp,Home      &     Furniture,Fossil
Watch,62.0,1.0,0.4,39.6,12.6,High,not_defined
2018-01-06,21:06:53,9.0,0.55476,Male,Mobile,New ,Auto & Accessories,Iron,,4.0,0.1,34.4,,e_wallet
2018-02-21,14:51:48,8.0,0.44850,Male,Web,Member,Electronic,Car      &     Bike
Care,192.0,,,162.5,1.9,Medium,money_order
</dataframe>
Update this initial code:
```python
TODO: import the required dependencies
import pandas as pd
Write code here
Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```
### QUERY
create a pie chart on customer login type column
Variable `dfs: list[pd.DataFrame]` is already declared.
At the end, declare "result" variable as a dictionary of type and value.
If you are asked to plot a chart, use "matplotlib" for charts, save as png.
Generate python code and return full updated code:
```

2025-03-14 22:11:55 [INFO] Code generated:

```

```python
TODO: import the required dependencies
import pandas as pd
import matplotlib.pyplot as plt
Assuming dfs is already defined and contains the DataFrame
For demonstration, let's use the first DataFrame in the list
df = dfs[0]
Count the occurrences of each customer login type
login_type_counts = df['Customer_Login_type'].value_counts()
Create a pie chart
plt.figure(figsize=(8, 6))
plt.pie(login_type_counts, labels=login_type_counts.index, autopct='%.1f%%', startangle=140)
plt.title('Customer Login Type Distribution')
plt.axis('equal') # Equal aspect ratio ensures that pie chart is circular.
Save the pie chart as a PNG file
plt.savefig('customer_login_type_distribution.png')
Declare result var:
result = {
 "type": "plot",
 "value": "customer_login_type_distribution.png"
}
```

```

```

2025-03-14 22:11:55 [INFO] Executing Step 4: CachePopulation
2025-03-14 22:11:55 [INFO] Executing Step 5: CodeCleaning
2025-03-14 22:11:55 [INFO] Saving charts to /Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png
2025-03-14 22:11:55 [INFO]
Code running:
```
df = dfs[0]
login_type_counts = df['Customer_Login_type'].value_counts()
plt.figure(figsize=(8, 6))
plt.pie(login_type_counts, labels=login_type_counts.index, autopct='%1.1f%%', startangle=140)
plt.title('Customer Login Type Distribution')
plt.axis('equal')
plt.savefig('/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png')
result = {'type': 'plot', 'value': '/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png'}
```
2025-03-14 22:11:55 [INFO] Executing Step 6: CodeExecution
2025-03-14 22:13:44 [INFO] Question: create a pir chart on customer login type column
2025-03-14 22:13:44 [INFO] Running PandasAI with openai LLM...
2025-03-14 22:13:44 [INFO] Prompt ID: 805e186c-08ac-427a-ae3e-dd6b7634146b
2025-03-14 22:13:44 [INFO] Executing Pipeline: GenerateChatPipeline
2025-03-14 22:13:44 [INFO] Executing Step 0: ValidatePipelineInput
2025-03-14 22:13:44 [INFO] Executing Step 1: CacheLookup
2025-03-14 22:13:44 [INFO] Executing Step 2: PromptGeneration
2025-03-14 22:13:44 [INFO] Using prompt: <dataframe>
dfs[0]:51290x16
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
2018-03-23,20:35:13,,89277,Male,Mobile,First SignUp,Electronic,Shirts,33.0,3.0,0.3,118.0,,Critical,money_order
2018-04-04,20:14:37,6.0,0,21457,Female,Mobile,Guest,Fashion,Car Mat,,4.0,0.1,120.6,12.8,,debit_card
2018-04-27,10:56:19,2.0,0.97589,Female,Web,New ,Auto & Accessories,Suits,117.0,,,104.8,15.6,Low,not_defined
</dataframe>
Update this initial code:
```
python
TODO: import the required dependencies
import pandas as pd
Write code here
Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```
### QUERY
create a pir chart on customer login type column
Variable `dfs: list[pd.DataFrame]` is already declared.
At the end, declare "result" variable as a dictionary of type and value.
If you are asked to plot a chart, use "matplotlib" for charts, save as png.
Generate python code and return full updated code:

```

2025-03-14 22:13:44 [INFO] Executing Step 3: CodeGenerator

2025-03-14 22:13:52 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 200 OK"

2025-03-14 22:13:52 [INFO] Prompt used:

```
<dataframe>
dfs[0]:51290x16
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
2018-03-23,20:35:13,,89277,Male,Mobile,First SignUp,Electronic,Shirts,33.0,3.0,0.3,118.0,,Critical,money_order
2018-04-04,20:14:37,6.0,21457,Female,Mobile,Guest,Fashion,Car Mat,,4.0,0.1,120.6,12.8,,debit_card
2018-04-27,10:56:19,2.0,97589,Female,Web,New ,Auto & Accessories,Suits,117.0,,,104.8,15.6,Low,not_defined
</dataframe>
```

Update this initial code:

```
'''python
# TODO: import the required dependencies
import pandas as pd
# Write code here
# Declare result var:
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
'''
```

QUERY

create a pie chart on customer login type column

Variable `dfs: list[pd.DataFrame]` is already declared.

At the end, declare "result" variable as a dictionary of type and value.

If you are asked to plot a chart, use "matplotlib" for charts, save as png.

Generate python code and return full updated code:

2025-03-14 22:13:52 [INFO] Code generated:

```
'''python
# TODO: import the required dependencies
import pandas as pd
import matplotlib.pyplot as plt
# Assuming dfs is already defined and contains the DataFrame
# For this example, we will use the first DataFrame in the list
df = dfs[0]
# Count the occurrences of each customer login type
login_type_counts = df['Customer_Login_type'].value_counts()
# Create a pie chart
plt.figure(figsize=(8, 6))
plt.pie(login_type_counts, labels=login_type_counts.index, autopct='%.1f%%', startangle=140)
plt.title('Customer Login Type Distribution')
plt.axis('equal') # Equal aspect ratio ensures that pie chart is circular.
# Save the pie chart as a PNG file
plt.savefig('customer_login_type_distribution.png')
# Declare result var:
result = {
    "type": "plot",
}
```

```

    "value": "customer_login_type_distribution.png"
}
```
2025-03-14 22:13:52 [INFO] Executing Step 4: CachePopulation
2025-03-14 22:13:52 [INFO] Executing Step 5: CodeCleaning
2025-03-14 22:13:52 [INFO] Saving charts to
/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png
2025-03-14 22:13:52 [INFO]
Code running:
```
df = dfs[0]
login_type_counts = df['Customer_Login_type'].value_counts()
plt.figure(figsize=(8, 6))
plt.pie(login_type_counts, labels=login_type_counts.index, autopct='%1.1f%%', startangle=140)
plt.title('Customer Login Type Distribution')
plt.axis('equal')
plt.savefig('/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png')
result = {'type': 'plot', 'value': '/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png'}
```
2025-03-14 22:13:52 [INFO] Executing Step 6: CodeExecution
2025-03-14 22:13:52 [INFO] Executing Step 7: ResultValidation
2025-03-14 22:13:52 [INFO] Answer: {'type': 'plot', 'value': '/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png'}
2025-03-14 22:13:52 [INFO] Executing Step 8: ResultParsing
2025-03-14 22:14:44 [INFO] Question: create a pie chart on payment methods column
2025-03-14 22:14:44 [INFO] Running PandasAI with openai LLM...
2025-03-14 22:14:44 [INFO] Prompt ID: 31158b43-3428-4e01-b5d5-c38ef15516cb
2025-03-14 22:14:44 [INFO] Executing Pipeline: GenerateChatPipeline
2025-03-14 22:14:44 [INFO] Executing Step 0: ValidatePipelineInput
2025-03-14 22:14:44 [INFO] Executing Step 1: CacheLookup
2025-03-14 22:14:44 [INFO] Executing Step 2: PromptGeneration
2025-03-14 22:14:44 [INFO] Using prompt: <dataframe>
dfs[0]:51290x16
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
2018-05-27,08:06:12,,46781,Female,Mobile,Member,Fashion,Beds,248.0.,0.1,10.8,4.9.,credit_card
2018-10-21,15:33:54,6.0,95273,Male,Web,New ,Home & Furniture,Jeans,220.0,1.0,0.3,32.0.,High,e_wallet
2018-11-30,14:32:51,4.0,33440,Female,Web,Guest,Electronic,Tablet,,4.0,,137.1,8.4,Critical,money_order
</dataframe>
Update this initial code:
```
python
# TODO: import the required dependencies
import pandas as pd
# Write code here
# Declare result var:

```

```
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```

```

### ### QUERY

create a pie chart on payment methods column

Variable `dfs: list[pd.DataFrame]` is already declared.

At the end, declare "result" variable as a dictionary of type and value.

If you are asked to plot a chart, use "matplotlib" for charts, save as png.

Generate python code and return full updated code:

```
2025-03-14 22:14:44 [INFO] Executing Step 3: CodeGenerator
```

```
2025-03-14 22:14:53 [INFO] HTTP Request: POST https://api.openai.com/v1/chat/completions "HTTP/1.1 200 OK"
```

```
2025-03-14 22:14:53 [INFO] Prompt used:
```

```
 <dataframe>
```

```
dfs[0]:51290x16
```

```
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
```

```
2018-05-27,08:06:12,46781,Female,Mobile,Member,Fashion,Beds,248.0,,0.1,10.8,4.9,,credit_card
```

```
2018-10-21,15:33:54,6.0,95273,Male,Web,New ,Home & Furniture,Jeans,220.0,1.0,0.3,32.0,,High,e_wallet
```

```
2018-11-30,14:32:51,4.0,33440,Female,Web,Guest,Electronic,Tablet,,4.0,,137.1,8.4,Critical,money_order
```

```
</dataframe>
```

Update this initial code:

```
```python
```

```
# TODO: import the required dependencies
```

```
import pandas as pd
```

```
# Write code here
```

```
# Declare result var:
```

```
type (possible values "string", "number", "dataframe", "plot"). Examples: { "type": "string", "value": f"The highest salary is {highest_salary}." } or { "type": "number", "value": 125 } or { "type": "dataframe", "value": pd.DataFrame({...}) } or { "type": "plot", "value": "temp_chart.png" }
```

```
```

```

### ### QUERY

create a pie chart on payment methods column

Variable `dfs: list[pd.DataFrame]` is already declared.

At the end, declare "result" variable as a dictionary of type and value.

If you are asked to plot a chart, use "matplotlib" for charts, save as png.

Generate python code and return full updated code:

```
2025-03-14 22:14:53 [INFO] Code generated:
```

```
```

```

```
# TODO: import the required dependencies
```

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
# Assuming dfs is already defined and contains the DataFrame
```

```
# For demonstration, let's create a sample DataFrame similar to the one provided
```

```
data = {
```

```
    "Order_Date": ["2018-05-27", "2018-10-21", "2018-11-30"],
```

```
    "Time": ["08:06:12", "15:33:54", "14:32:51"],
```

```

    "Aging": [None, 6.0, 4.0],
    "Customer_Id": [46781, 95273, 33440],
    "Gender": ["Female", "Male", "Female"],
    "Device_Type": ["Mobile", "Web", "Web"],
    "Customer_Login_type": ["Member", "New", "Guest"],
    "Product_Category": ["Fashion", "Home & Furniture", "Electronic"],
    "Product": ["Beds", "Jeans", "Tablet"],
    "Sales": [248.0, 220.0, None],
    "Quantity": [None, 1.0, 4.0],
    "Discount": [0.1, 0.3, None],
    "Profit": [10.8, 32.0, 137.1],
    "Shipping_Cost": [4.9, None, 8.4],
    "Order_Priority": [None, "High", "Critical"],
    "Payment_method": ["credit_card", "e_wallet", "money_order"]
}
# Create DataFrame
dfs = [pd.DataFrame(data)]
# Extract the Payment_method column and count occurrences
payment_counts = dfs[0]['Payment_method'].value_counts()
# Create a pie chart
plt.figure(figsize=(8, 6))
plt.pie(payment_counts, labels=payment_counts.index, autopct='%1.1f%%', startangle=140)
plt.title('Distribution of Payment Methods')
plt.axis('equal') # Equal aspect ratio ensures that pie chart is circular.
plt.savefig('payment_methods_distribution.png') # Save the pie chart as a PNG file
plt.close() # Close the plot to free up memory
# Declare result var:
result = {
    "type": "plot",
    "value": "payment_methods_distribution.png"
}
```

```

2025-03-14 22:14:53 [INFO] Executing Step 4: CachePopulation  
2025-03-14 22:14:53 [INFO] Executing Step 5: CodeCleaning

|            |          |        |        |        |    |
|------------|----------|--------|--------|--------|----|
| 2025-03-14 | 22:14:53 | [INFO] | Saving | charts | to |
|------------|----------|--------|--------|--------|----|

/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp\_chart.png

2025-03-14 22:14:53 [INFO]  
Code running:  
```
data = {'Order_Date': ['2018-05-27', '2018-10-21', '2018-11-30'], 'Time': ['08:06:12', '15:33:54', '14:32:51'], 'Aging': [None, 6.0, 4.0], 'Customer_Id': [46781, 95273, 33440], 'Gender': ['Female', 'Male', 'Female'], 'Device_Type': ['Mobile', 'Web', 'Web'], 'Customer_Login_type': ['Member', 'New', 'Guest'], 'Product_Category': ['Fashion', 'Home & Furniture', 'Electronic'], 'Product': ['Beds', 'Jeans', 'Tablet'], 'Sales': [248.0, 220.0, None], 'Quantity': [None, 1.0, 4.0], 'Discount': [0.1, 0.3, None], 'Profit': [10.8, 32.0, 137.1], 'Shipping_Cost': [4.9, None, 8.4], 'Order_Priority': [None, 'High', 'Critical'], 'Payment_method': ['credit_card', 'e_wallet', 'money_order']}

payment_counts = dfs[0]['Payment_method'].value_counts()
plt.figure(figsize=(8, 6))
plt.pie(payment_counts, labels=payment_counts.index, autopct='%1.1f%%', startangle=140)

```

plt.title('Distribution of Payment Methods')
plt.axis('equal')
plt.savefig('/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png')
plt.close()
result = {'type': 'plot', 'value': '/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png'}
```
2025-03-14 22:14:53 [INFO] Executing Step 6: CodeExecution
2025-03-14 22:14:53 [INFO] Executing Step 7: ResultValidation
2025-03-14 22:14:53 [INFO] Answer: {'type': 'plot', 'value': '/Users/jsujanchowdary/Downloads/langchain-ask-csv-main/exports/charts/temp_chart.png'}
2025-03-14 22:14:53 [INFO] Executing Step 8: ResultParsing
2025-03-14 22:15:34 [INFO] Question: create a bar graph chart on payment methods column
2025-03-14 22:15:34 [INFO] Running PandasAI with openai LLM...
2025-03-14 22:15:34 [INFO] Prompt ID: d5e7e2ab-77b3-4b54-8abd-e7fdd0e1979a
2025-03-14 22:15:34 [INFO] Executing Pipeline: GenerateChatPipeline
2025-03-14 22:15:34 [INFO] Executing Step 0: ValidatePipelineInput
2025-03-14 22:15:34 [INFO] Executing Step 1: CacheLookup
2025-03-14 22:15:34 [INFO] Executing Step 2: PromptGeneration
2025-03-14 22:15:34 [INFO] Using prompt: <dataframe>
dfs[0]:51290x16
Order_Date,Time,Aging,Customer_Id,Gender,Device_Type,Customer_Login_type,Product_Category,Product,Sales,
Quantity,Discount,Profit,Shipping_Cost,Order_Priority,Payment_method
2018-03-05,18:44:54,10.0,90751,Male,Mobile,New ,Auto & Accessories,Suits,140.0,,,34.0,4.4,,debit_card
2018-04-28,23:44:42,4.0,87775,Male,Mobile,Member,Home & Furniture,Iron,,1.0,0.5,54.0,5.5,High,credit_card
2018-10-22,18:34:05,,37959,Female,Web,First SignUp,Fashion,Titak watch,54.0,5.0,0.4,35.8,,Low,money_order
</dataframe>

```