

AIRBNB DATA CLEANING, PROCESSING, AND ADVANCED ANALYSIS

Shubham Raj – G23AI2028

Bhavesh Arora – G23AI2126

Niraj Bhagchandani – G23AI2087

Jai Singh Khushvah – G23AI2018

Paras Panda – G23AI2117

Jatin Shrivastava – G23AI2094

CAPSTONE PROJECT

This Report Presented for the PG Diploma in Data
Engineering

Department of Computer Science IIT Jodhpur



CONTENTS

- 1 Project Background**

- 2 Project Approach**

- 3 Key Deliverables**

- 4 Comprehensive Overview
of Airbnb Datasets**

- 5 Different Phases of Project**

- 6 Screenshots**

- 7 Future work**
- Recommendation &Conclusion**

PROJECT BACKGROUND



Onboarding and Initial Data Handling

1

Understand, clean, and augment the Airbnb dataset, and perform exploratory data analysis (EDA)



Data Ingestion and Optimization

2

Develop efficient data ingestion pipelines and optimize data storage.



Data Transformation

3

Transform data to meet business requirements and enable advanced analytics



Data Warehousing, Reporting, and Visualization

4

Understand, clean, and augment the Airbnb dataset, and perform exploratory data analysis (EDA)

KEY DELIVERABLES



01

Monthly Progress Reports : Detailed reports outlining progress on each problem statement, including methodologies, challenges, and preliminary results.



02

Final Data Warehouse : A fully populated and optimized data warehouse schema.



03

Predictive Models : Trained and validated machine learning models with performance metrics.



04

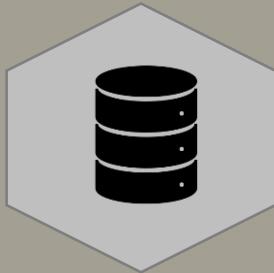
BI Dashboards: Interactive dashboards for exploring data and insights.



05

Comprehensive Final Report : A detailed report summarizing the entire project, including methodologies, findings, and recommendations.

COMPREHENSIVE OVERVIEW OF AIRBNB DATASETS

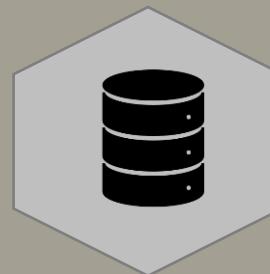


Calendar.csv

- Number of Rows: 1,393,570
- Number of Columns: 7
- Key Column Names:
 - Listing_id,
 - Date,
 - Available
 - Price
 - Minimum nights
 - Maximum nights
 - Available units



Calender

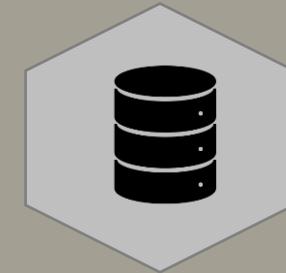


Listings.csv

- Number of Rows: 3,818
- Number of Columns: 92
- Key Column Names:
 - id
 - listing_url
 - scrape_id,
 - Last_scraped
 - Name
 - host_id
 - host_name,



Listing



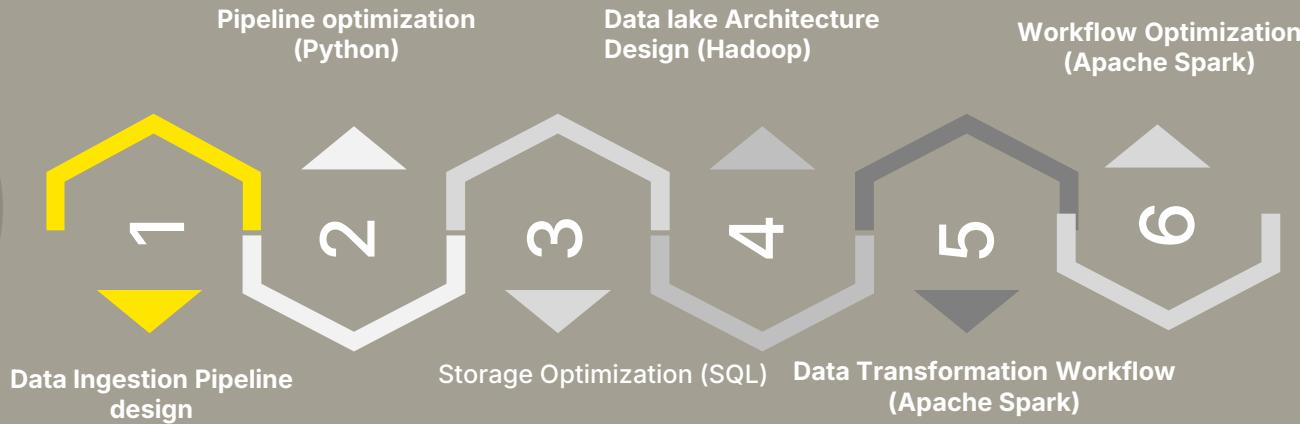
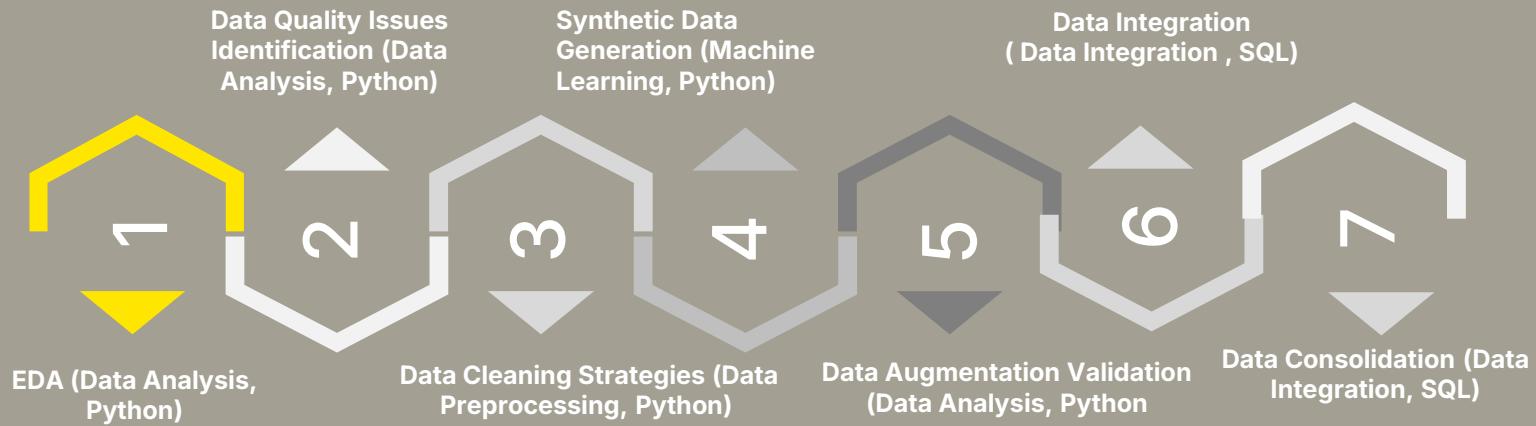
Reviews.csv

- Number of Rows: 84,849
- Number of Columns: 7
- Key Column Names:
 - listing_id
 - id
 - Date
 - reviewer_id
 - reviewer_name
 - Comments
 - sentiment score

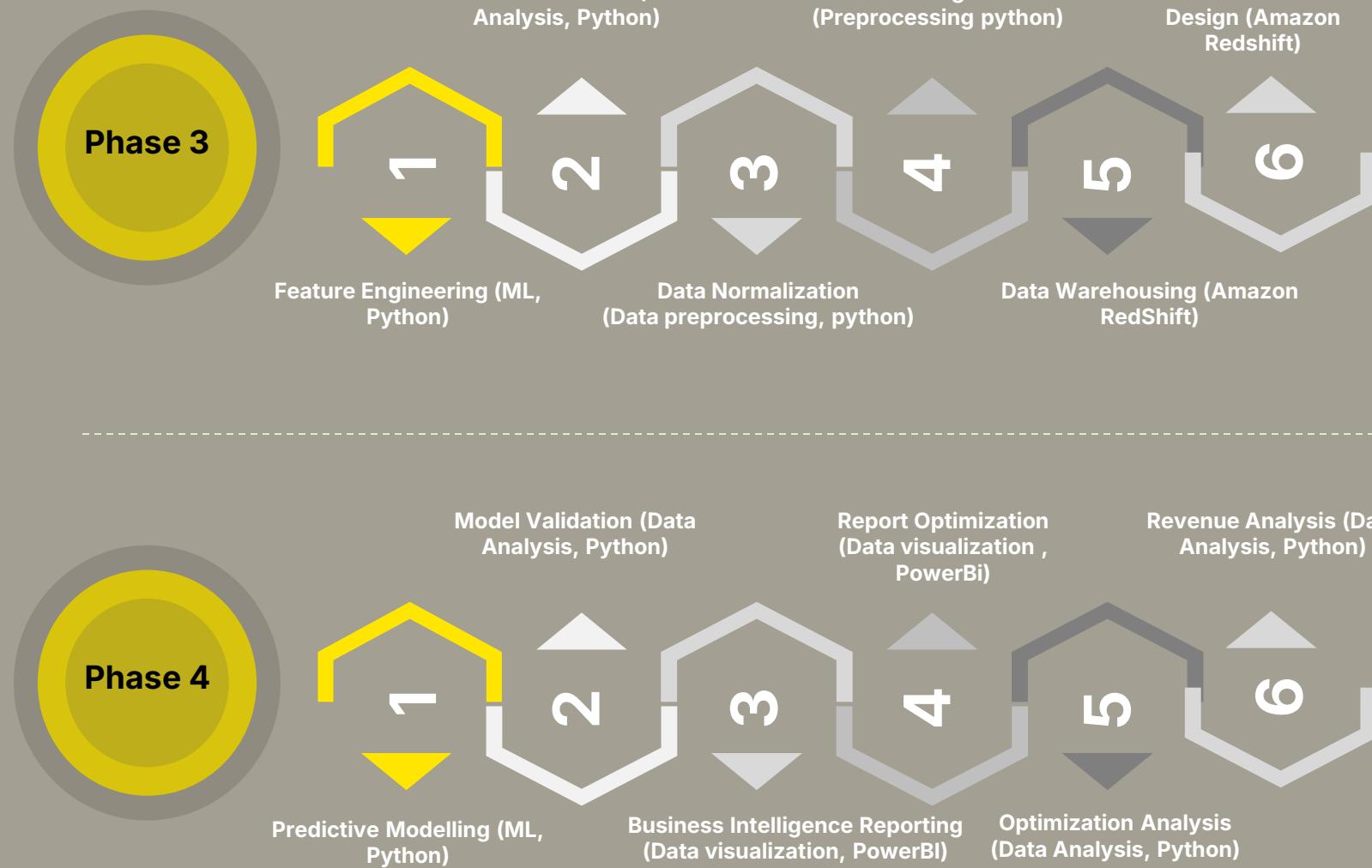


Review

DIFFERENT PHASES OF PROJECT



DIFFERENT PHASES OF PROJECT



Some Useful Links

GitHub Links:

<https://github.com/shubham14p3/IITJ-AIRBNB-DATA-CLEANING-PROCESSING-AND-ADVANCED-ANALYSIS-CAPSTONE-PROJECT/tree/main/backend>

Front End Link:

<http://3.222.77.245:5173/>

Back End Link:

<http://44.198.64.142:5000/>



Phase

01

Onboarding and Initial
Data Handling

PROBLEM STATEMENT & PROPOSED SOLUTIONS

1



Comprehensive Data Profiling and Cleaning

Conduct EDA on calendar.csv, listings.csv, and reviews.csv. Identify and visualize data quality issues.

2



Data Augmentation and Synthesis

Develop techniques to augment the dataset by creating synthetic data.

3



Data Integration and Consolidation

Integrate calendar.csv, listings.csv, and reviews.csv into a consolidated dataset.

Instructions

- Perform EDA to identify missing values, outliers, and anomalies using visualizations like heatmaps and box plots.
- Implement data cleaning strategies, including imputing missing values, treating outliers, and standardizing data formats.

- Develop machine learning models (e.g., GANs or SMOTE) to predict missing values and generate synthetic data points.
- Validate the augmented data using statistical methods to ensure robustness.

- Merge the datasets, resolve inconsistencies such as duplicate records, and remove redundant information.
- Optimize the consolidated dataset for subsequent analysis by standardizing data types and normalizing values.

Solution

- **EDA Insights** : Analysed datasets with summaries, heatmaps for missing data, and box plots for outliers.
- **Key Findings** : Missing values in critical columns like price and review_scores. Outliers in price and availability_365.
- **Data Cleaning** : Imputed missing values (median for numeric, mode for categorical). Capped outliers at the 95th percentile. Standardized date and currency formats.
- **Documentation** : Captured key processes, code snapshots, and a summary table for cleaned data.

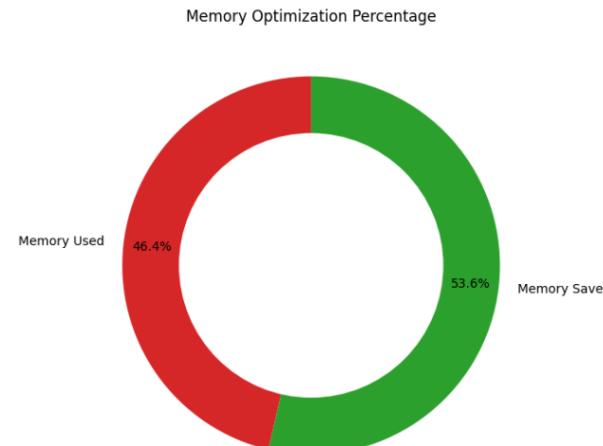
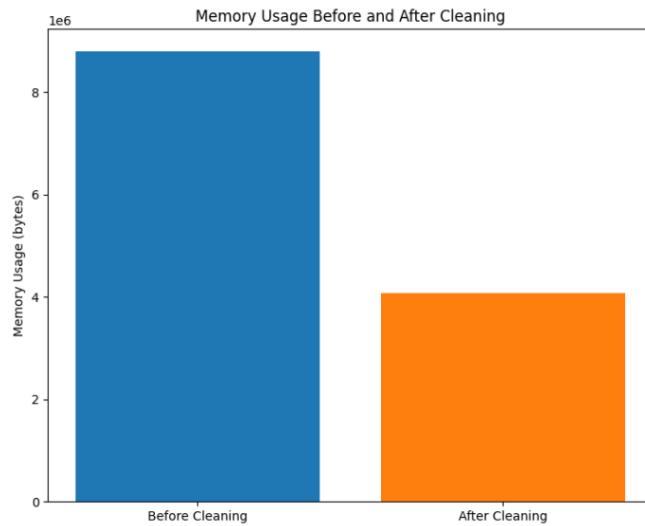
- **Data Augmentation** : Used SMOTE to address class imbalance in categorical variables and GANs to synthesize data for continuous columns, mitigating the impact of missing values.
- **Validation** : Conducted statistical tests (e.g., Chi-square and t-tests) to verify the synthetic data's distribution and cross-validated to ensure alignment with original patterns.
- **Outcome** : Successfully filled gaps in listings.csv and reviews.csv with synthetic data, enhancing dataset robustness for analysis.

- **Merging Process** : Identified primary keys (e.g., listing_id), resolved duplicates, and addressed data type discrepancies during dataset merging.
- **Data Standardization** : Ensured consistent formats for dates and numeric columns and normalized values to align with expected ranges for smooth analysis.

Listings Dataset

```
Changes in Missing Values:  
      Missing Values Percentage  
id           0    0.0  
listing_url   0    0.0  
name          0    0.0  
host_id        0    0.0  
host_name      0    0.0  
host_since     0    0.0  
neighbourhood 0    0.0  
latitude       0    0.0  
longitude      0    0.0  
number_of_reviews 0    0.0  
review_scores_rating 0    0.0  
instant_bookable 0    0.0  
room_type_id    0    0.0  
cancellation_policy_id 0    0.0
```

```
Memory Usage Optimization:  
Memory reduced by: 4718414 bytes  
Memory optimized by: 53.64%
```



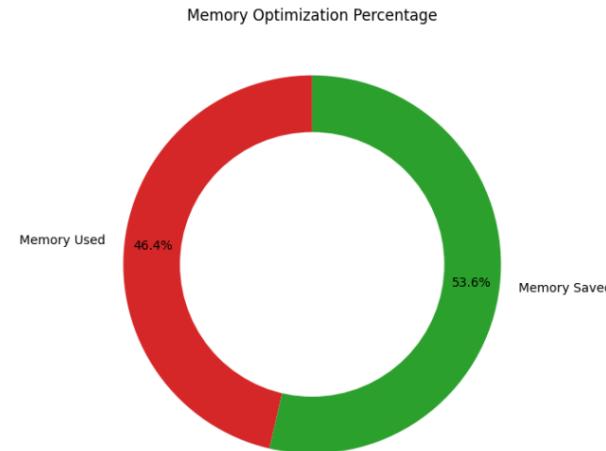
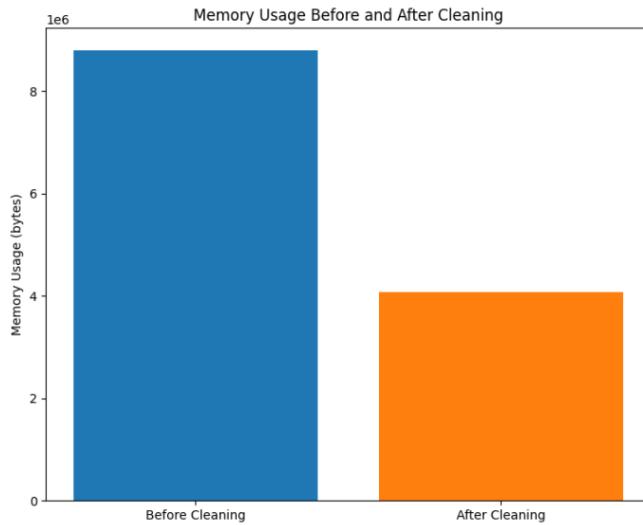
53% Memory
Saved

Reviews Dataset

Changes in Missing Values:

	Missing Values	Percentage
id	0	0.0
listing_url	0	0.0
name	0	0.0
host_id	0	0.0
host_name	0	0.0
host_since	0	0.0
neighbourhood	0	0.0
latitude	0	0.0
longitude	0	0.0
number_of_reviews	0	0.0
review_scores_rating	0	0.0
instant_bookable	0	0.0
room_type_id	0	0.0
cancellation_policy_id	0	0.0

Memory Usage Optimization:
Memory reduced by: 4718414 bytes
Memory optimized by: 53.64%



53% Memory
Saved

Calendar Dataset

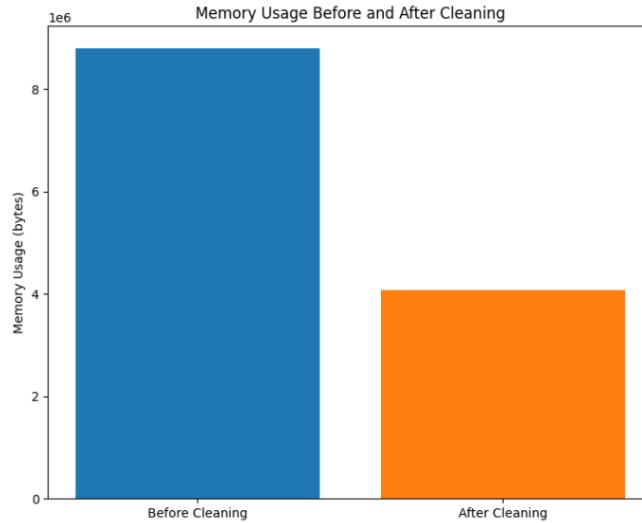
Changes in Missing Values:

	Missing Values	Percentage
id	0	0.0
listing_url	0	0.0
name	0	0.0
host_id	0	0.0
host_name	0	0.0
host_since	0	0.0
neighbourhood	0	0.0
latitude	0	0.0
longitude	0	0.0
number_of_reviews	0	0.0
review_scores_rating	0	0.0
instant_bookable	0	0.0
room_type_id	0	0.0
cancellation_policy_id	0	0.0

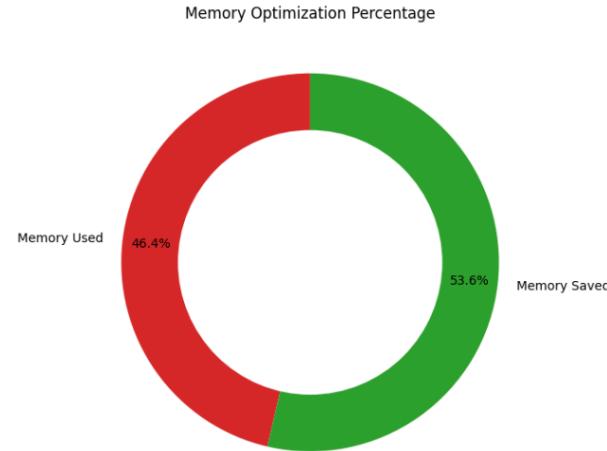
Memory Usage Optimization:

Memory reduced by: 4718414 bytes

Memory optimized by: 53.64%



Memory Optimization Percentage



53% Memory
Saved

Phase

02

Data Ingestion and
Optimization

PROBLEM STATEMENT & PROPOSED SOLUTIONS

4



Scalable Data Ingestion Pipeline

Design and implement a scalable data ingestion pipeline using Python, Pandas, and SQL.

5



Advanced Storage Optimization Techniques

Implement storage optimization techniques such as data partitioning, indexing, and compression.

Instructions

- Incorporate error handling, logging, and monitoring features for reliability.
- Optimize the pipeline for performance, including batch processing and real-time data ingestion.

- Apply storage optimization in a cloud-based data storage solution like AWS RDS or Google Big Query.

Solution

- **Pipeline Development:** Designed a scalable data ingestion pipeline using Python, Pandas, and SQL for efficient data handling and storage.
- **Error Handling and Logging:** Incorporated robust error handling and logging mechanisms to enhance reliability and troubleshoot issues effectively.
- **Performance Optimization:** Implemented batch processing for large datasets and real-time ingestion features to support dynamic data workflows.
- **Monitoring and Alerts:** Enabled proactive monitoring with dashboards or alerts to track pipeline health and ensure smooth operation.

- **Data Partitioning:** Implemented partitioning techniques in cloud-based storage like AWS RDS or Google BigQuery to organize data for faster query performance.
- **Indexing:** Applied indexing on critical fields to enhance data retrieval speeds and support high-performance queries.
- **Data Compression:** Enabled compression for large datasets to reduce storage costs while maintaining data integrity and accessibility.
- **Cloud-Based Optimization:** Leveraged tools in AWS RDS or BigQuery for seamless integration of these optimizations within a scalable cloud environment.

Phase

03

Data Transformation

PROBLEM STATEMENT & PROPOSED SOLUTIONS

6



Complex Data Transformation Workflows

Design data transformation workflows using tools like Apache Spark or Apache Airflow.

- Handle data aggregation, filtering, and enrichment processes to prepare for advanced analytics.
- Ensure the workflows are scalable and efficient.

7



Feature Engineering for Machine Learning

Develop techniques for feature engineering to extract meaningful features from the dataset.

Instructions

- Extract temporal features from calendar.csv and textual features from reviews.csv.
- Validate the effectiveness of these features for machine learning models.

8



Advanced Data Normalization and Encoding

Implement data normalization techniques for machine learning model preparation.

- Apply normalization techniques such as min-max scaling and Z-score normalization.
- Use encoding methods like one-hot encoding, label encoding, and embeddings for categorical variables.

Solution

- Data Aggregation:** Use Spark's groupBy and agg functions to compute metrics (e.g., total sales, average values).
- Data Filtering:** Eliminate unwanted or erroneous records using the filter function
- Partitioning and Caching:** Ensure datasets are partitioned by frequently used keys (e.g., date or region).
Caching: Cache intermediate results to reduce redundant computations.
- Validation:** Perform checks to ensure data accuracy and completeness using test cases or by comparing against

- Feature Engineering :** Create and preprocess new features, encode categories, and handle missing data.
- Model Training & Evaluation :** Train models and assess performance using metrics like accuracy and precision.
- Impact Analysis :** Compare model performance before and after feature engineering and identify key features driving improvements.

- Normalization & Encoding:** Apply Scikit-Learn's Standard Scaler and OneHotEncoder to scale numerical features and encode categorical variables.
- Model Performance:** Train models using the preprocessed data and measure key metrics like accuracy, precision, or recall.

Phase

04

**Data Warehousing,
Reporting, and
Visualization**

PROBLEM STATEMENT & PROPOSED SOLUTIONS

9



Enterprise Data Warehousing

Design an enterprise-level data warehouse schema using a cloud-based solution.

10



Advanced Business Intelligence Reporting

Create advanced business intelligence reports and dashboards using tools like Power BI or Tableau.

Instructions

- Create star and snowflake schemas to support complex querying and analytics.
- Focus on key business metrics such as occupancy rates, revenue per listing, and customer satisfaction..

Solution

- **Schema Design:** Create an efficient data warehouse schema using Redshift or BigQuery.
- **ETL Process:** Implement ETL pipelines for data extraction, transformation, and loading.
- **Quality & Performance:** Ensure data quality and optimize query performance within the data warehouse.

- **Dashboard Creation:** Build interactive dashboards using tools like Tableau or Power BI.
- **Data Integration:** Integrate and aggregate data for real-time reporting.
- **Insights & Optimization:** Derive actionable insights and optimize visualizations for clarity.

PROBLEM STATEMENT & PROPOSED SOLUTIONS

11



Predictive Analytics and Machine Learning

Develop predictive models using machine learning algorithms.

12



Optimization and Revenue Analysis

Perform advanced data analysis to identify optimization areas and potential revenue growth.

Instructions

- Use regression, classification, and clustering techniques to derive insights from the data.

- Use techniques like A/B testing, cohort analysis, and time series analysis to identify optimization areas.

Solution

- **Model Building:** Develop and train predictive models using appropriate algorithms.
- **Model Validation:** Validate model performance through testing and cross-validation.
- **Iterative Improvement:** Continuously refine the model based on performance metrics and document the development process and results.

- **Data Analysis:** Analyze data to identify trends, patterns, and key insights.
- **Actionable Recommendations:** Provide recommendations for optimization and revenue growth based on analysis.
- **Documentation:** Document findings and the impact of recommendations for future reference.

SCREENSHOTS

- ❖ JIRA Project Allocation and Tracking
- ❖ Github project Upload and Contribution
- ❖ Screenshots from Google cloud and AWS (RDS)

JIRA PROJECT ALLOCATION AND TRACKING

Projects / Airbnb-IITJ-Capstone-Project-Group-20
SCRUM Sprint 2
Capstone Project 2024 IITJ

11 days | ⚡ ⭐ 🔍 ⏪ Complete sprint ⏪

PLANNING

- Summary **NEW**
- Timeline
- Backlog
- Board**
- Forms
- Goals
- + Add view

DEVELOPMENT

- Code

Project pages

TO DO 2

- 11. Problem Statement 11: Predictive Analytics and Machine Learning (SCRUM-38)
- 12. Problem Statement 12: Optimization and Revenue Analysis (SCRUM-39)

+ Create

IN PROGRESS 9

- 3. Problem Statement 3: Data Integration and Consolidation (SCRUM-30)
- 4. Problem Statement 4: Scalable Data Ingestion Pipeline (SCRUM-31)
- 6. Problem Statement 6: Complex Data Transformation Workflows (SCRUM-32)

COMPLETED 2

- 2. Problem Statement 2: Data Augmentation and Synthesis (SCRUM-29)
- 5. Problem Statement 5: Advanced Storage Optimization Techniques (SCRUM-32)

DONE 1

- 1. Problem Statement 1: Comprehensive Data Cleaning (SCRUM-28)

GROUP BY None ⏪ ⏪

Quickstart

Your work ▾ Projects ▾ Filters ▾ Dashboards ▾ Teams ▾ More ▾ Create Upgrade Search 4 ⓘ ⓘ ⓘ ⓘ

Airbnb-IITJ-Capstone-Pr... Software project

To customize user access, such as roles and permissions, upgrade your plan to Standard.

Upgrade Learn more about managing access

Back to project

Details

Access

Notifications

Automation

Issue types

Features

Board

Toolchain

Apps

Search roles Roles

Name	Email	Role	Action
BA Bhavesh Arora(G23AI2126)	g23ai2126@iitj.ac.in	Administrator	⋮
G g23ai2087	-	Administrator	⋮
JK Jai Singh Kushwah(G23AI2018)	g23ai2018@iitj.ac.in	Administrator	⋮
I Jatin Shrivastava (G23AI2094)	g23ai2094@iitj.ac.in	Administrator	⋮
P Paras Panda (G23AI2117)	g23ai2117@iitj.ac.in	Administrator	⋮
S Shubham Raj (G23AI2028)	g23ai2028@iitj.ac.in	Administrator	⋮

Quickstart

Jira URL:
<https://iitj-team-data-crunchers.atlassian.net/jira/software/projects/SCRUM/boards/1>

GITHUB PROJECT UPLOAD AND CONTRIBUTION

IITJ-AIRBNB-DATA-CLEANING-PROCESSING-AND-ADVANCED-ANALYSIS-...

Public Watch 1 Fork 0 Star 0

main 1 Branch 0 Tags Go to file + <> Code

shubham14p3 Adding search schema for uiniq ways 4bb3342 · 9 hours ago 27 Commits

Documents Adding FOuter 2 months ago

backend Adding Uniqueness to data 10 hours ago

public Adding UI 2 months ago

src Adding search schema for uiniq ways 9 hours ago

.gitignore Adding Files for backend 20 hours ago

README.md Adding FOuter 2 months ago

eslint.config.js Adding UI 2 months ago

index.html Adding UI 2 months ago

package-lock.json Adding search schema for uiniq ways 9 hours ago

package.json Adding search schema for uiniq ways 9 hours ago

About

IITJ AIRBNB DATA CLEANING, PROCESSING AND ADVANCED ANALYSIS

Readme Activity 0 stars 1 watching 0 forks Report repository

Releases

No releases published Create a new release

Packages

No packages published Publish your first package

README

Bhagchandani Niraj

- Github: [@BhagchandaniNiraj](#)
- Linkedin: [Niraj Bhagchandani](#)

Bhavesh Arora

- Github: [@BhaveshArora](#)
- Linkedin: [Bhavesh Arora](#)

Jai Singh Kushwah

- Github: [@JaiSinghKushwah](#)
- Linkedin: [Jai Singh Kushwah](#)

PARAS PANDA

- Github: [@PARASPANDA](#)
- Linkedin: [PARAS PANDA](#)

**JATIN SHRIVAS **

- Github: [@JATINSHRIVAS](#)
- Linkedin: [JATIN SHRIVAS](#)

GITHUB URL:
<https://github.com/shubham14p3/IITJ-AIRBNB-DATA-CLEANING-PROCESSING-AND-ADVANCED-ANALYSIS-CAPSTONE-PROJECT>

Apps

Atlassian Marketplace

Manage apps

App requests

Promotions

OAuth credentials

Apps

 GitHub Connector

GitHub for Jira

Configure

GitHub configuration

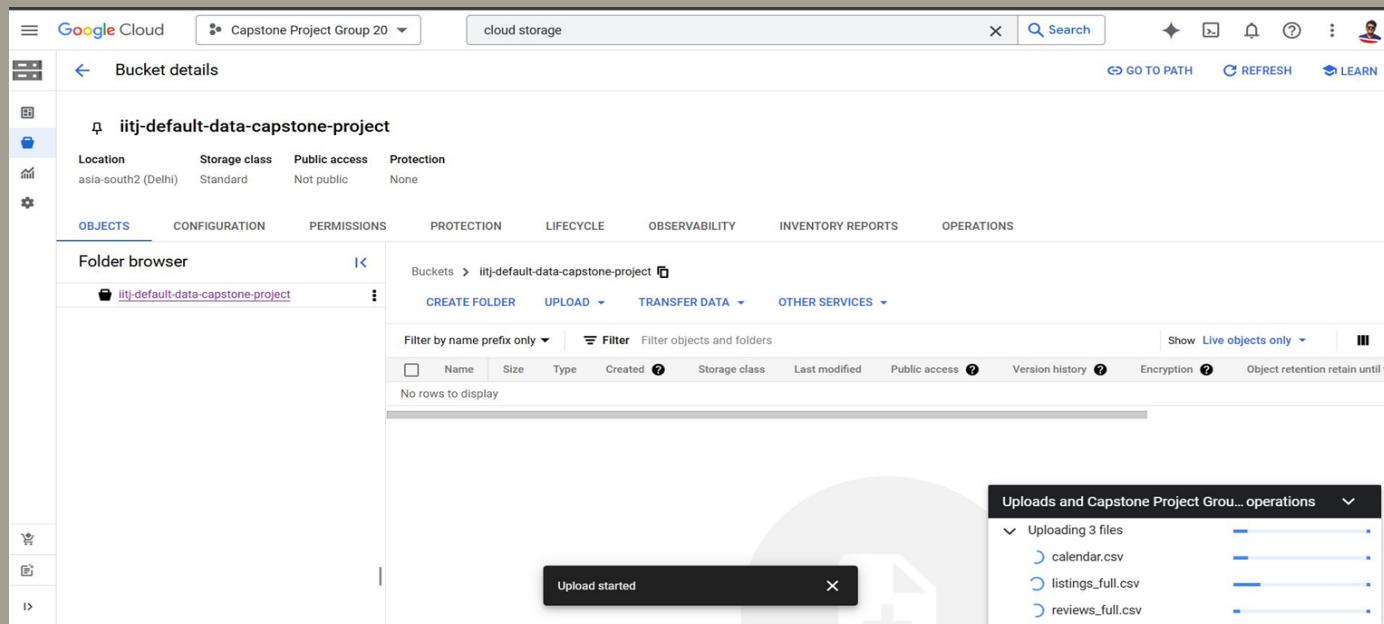
Connect a GitHub organization ▾

Connecting GitHub to Jira allows you to view development activity in the context of your Jira project and issues. To send development data from GitHub to Jira, your team must include issue keys in branch names, commit messages, and pull request titles.

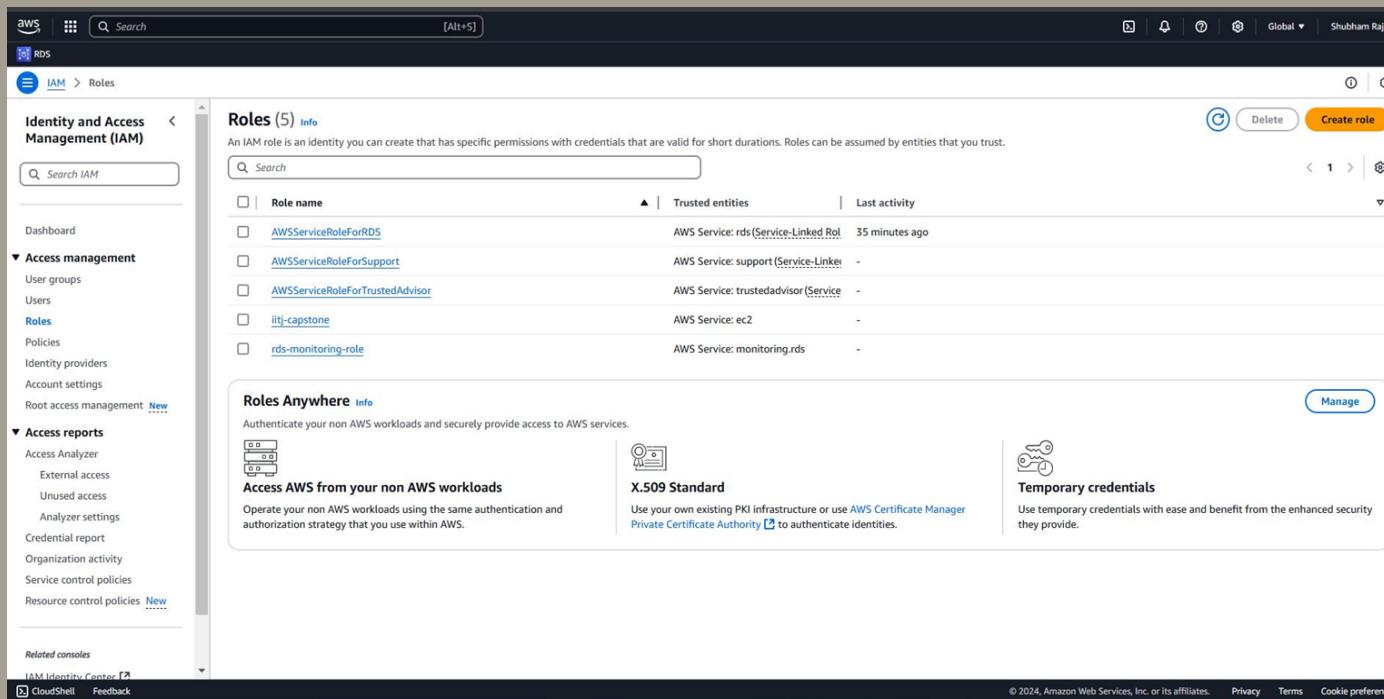
Even if your organization is still backfilling historical data, you can start using issue keys in your development work immediately.

Connected organization	Repository access	Backfill status	Permissions	Settings
 shubham14p3	Only select repos 1 	FINISHED Backfilled from: 6/9/24 	FULL ACCESS	...

SCREENSHOTS FROM GOOGLE CLOUD AND AWS (RDS)



The screenshot shows the Google Cloud Bucket details page for the bucket "iitj-default-data-capstone-project". The bucket is located in "asia-south2 (Delhi)", has a "Standard" storage class, and "Not public" public access. Protection is set to "None". The "OBJECTS" tab is selected, showing a "Folder browser" for the root folder. A modal window displays "Upload started" over the list of files being uploaded: "calendar.csv", "listings_full.csv", and "reviews_full.csv".



The screenshot shows the AWS IAM Roles page. The left sidebar shows "Identity and Access Management (IAM)" with "Roles" selected. The main content area shows a list of roles: "AWSServiceRoleForRDS", "AWSServiceRoleForSupport", "AWSServiceRoleForTrustedAdvisor", "iitj-capstone", and "rds-monitoring-role". Below this, there are sections for "Roles Anywhere" (info), "Access AWS from your non AWS workloads" (info), and "Temporary credentials" (info). The bottom of the page includes standard AWS navigation links like CloudShell and Feedback.

SCREENSHOTS FROM GOOGLE CLOUD AND AWS (RDS)

The screenshot shows the AWS IAM User Details page for the user 'shubham14p3'. The left sidebar includes links for Identity and Access Management (IAM), Access management (Users, Roles, Policies, Identity providers, Account settings, Root access management), Access reports (Access Analyzer, External access, Unused access, Analyzer settings, Credential report, Organization activity, Service control policies, Resource control policies), and Related consoles (IAM, Identity Center). The main content area displays the user's ARN (arn:aws:iam::329599615910:user/shubham14p3), which is disabled for console access. It also shows the user was created on December 06, 2024, at 21:10 UTC+05:30, with no last console sign-in. Two access keys are listed: 'Access key 1' (AKIAUZPNLAOTANJGJYHXJ - Active, Never used, Created today) and 'Access key 2' (Create access key). The 'Permissions' tab is selected, showing four attached policies: AmazonRDSFullAccess, AmazonS3FullAccess, CloudWatchFullAccess, and CloudWatchFullAccessV2, all of which are AWS managed and attached directly.

The screenshot shows the AWS S3 Upload status page. A large blue progress bar indicates an upload is in progress, with the message: 'Uploading', 'Total remaining: 3 files: 310.0 MB (99.77%)', 'Estimated time remaining: an hour', and 'Transfer rate: 79.6 kB/s'. Below this, a 'Cancel' button is visible. The main content area is titled 'Upload: status' and includes a summary table:

Destination	Succeeded	Failed
s3://litj-data-ingestion-bucket	0 files, 720.0 KB (0.23%)	0 files, 0 B (0%)

The 'Files and folders' tab is selected, showing a list of three files: 'calendar.csv', 'listings_full.csv', and 'reviews_full.csv', each with a size of approximately 171.7 MB, 45.8 MB, and 93.2 MB respectively. The status for all three files is 'Pending'.

The screenshot shows the AWS RDS Databases page. On the left, there's a sidebar with options like Dashboard, Databases (which is selected), Query Editor, Performance insights, Snapshots, Exports in Amazon S3, Automated backups, and Reserved instances. The main area has a title 'Databases (1)'. A search bar says 'Filter by databases'. Below it is a table with columns: DB identifier, Status, Role, Engine, Region ..., Size, Recommendations, CPU, Current ..., and Maintenance. There's one row for a database named 'itj' with status 'Creating', engine 'MySQL Co...', region 'us-east-1d', size 'db.t4g.micro', and maintenance 'none'. At the top right, there are buttons for 'Group resources', 'Modify', 'Actions', 'Restore from S3', and 'Create database'.

This screenshot is similar to the first one, showing the 'Databases (1)' page. However, there is a prominent green success message box at the top that says 'Successfully created database itj'. It also includes a note: 'You can use settings from itj to simplify configuration of suggested database add-ons while we finish creating your DB for you.' Below the message is the same table as the first screenshot, showing the database 'itj' now in the 'Available' state. The rest of the interface is identical to the first screenshot.

This screenshot is identical to the first one, showing the 'Databases (1)' page. The database 'itj' is still listed with the 'Creating' status. The rest of the interface, including the sidebar and the table structure, is the same as the other screenshots.

Uploading 0%

Total remaining: 3 files: 310.0 MB (99.77%)
Estimated time remaining: an hour
Transfer rate: 79.6 KB/s

Cancel Close

Upload: status

ⓘ After you navigate away from this page, the following information is no longer available.

Summary

Destination	Succeeded	Failed
s3://liltj-data-ingestion-bucket	0 files, 720.0 KB (0.23%)	0 files, 0 B (0%)

Files and folders Configuration

Files and folders (3 total, 310.7 MB)

Find by name

Name	Folder	Type	Size	Status	Error
calendar.csv	-	text/csv	171.7 MB	In progress (0%)	-
listings_full.csv	-	text/csv	45.8 MB	Pending	-
reviews_full.csv	-	text/csv	93.2 MB	Pending	-

SCREENSHOTS FROM GOOGLE CLOUD AND AWS (RDS)

RDS > Create database

MySQL

Additional configuration

Database options

Initial database name [Info](#) iitj

If you do not specify a database name, Amazon RDS does not create a database.

DB parameter group [Info](#) default.mysql8.0

Option group [Info](#) default:mysql-8-0

Backup

Enable automated backups

Creates a point-in-time snapshot of your database

Encryption

Enable encryption

Choose to encrypt the given instance. Master key IDs and aliases appear in the list after they have been created using the AWS Key Management Service console. [Info](#)

Log exports

Select the log types to publish to Amazon CloudWatch Logs

Audit log
 Error log
 General log
 Slow query log

IAM role

The following service-linked role is used for publishing logs to CloudWatch Logs.

RDS service-linked role

CloudShell Feedback

© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

MySQL

MySQL is the most popular open source database in the world. MySQL on RDS offers the rich features of the MySQL community edition with the flexibility to easily scale compute resources or storage capacity for your database.

- Supports database size up to 64 TiB.
- Supports General Purpose, Memory Optimized, and Burstable Performance instance classes.
- Supports automated backup and point-in-time recovery.
- Supports up to 15 Read Replicas per instance, within a single Region or 5 read replicas cross-region.



>



>

▼ Additional configuration

Database options, encryption turned off, backup turned off, backtrack turned off, maintenance, CloudWatch Logs, delete protection turned off.

Database options

Initial database name [Info](#)

 iitj

If you do not specify a database name, Amazon RDS does not create a database.

DB parameter group [Info](#)

 default.mysql8.0

Option group [Info](#)

 default:mysql-8.0

Backup

Enable automated backups

Creates a point-in-time snapshot of your database

Encryption

Enable encryption

Choose to encrypt the given instance. Master key IDs and aliases appear in the list after they have been created using the AWS Key Management Service console. [Info](#)

Log exports

Select the log types to publish to Amazon CloudWatch Logs

- Audit log
- Error log
- General log
- Slow query log

IAM role

The following service-linked role is used for publishing logs to CloudWatch Logs.

RDS service-linked role

MySQL

MySQL is the most popular open source database in the world. MySQL on RDS offers the rich features of the MySQL community edition with the flexibility to easily scale compute resources or storage capacity for your database.

- Supports database size up to 64 TiB.
- Supports General Purpose, Memory Optimized, and Burstable Performance instance classes.
- Supports automated backup and point-in-time recovery.
- Supports up to 15 Read Replicas per instance, within a single Region or 5 read replicas cross-region.



us-east-1.console.aws.amazon.com/s3/buckets/iitj-data-ingestion-buckets?region=us-east-1&bucketType=general&tab=objects

aws RDS Amazon S3 Buckets iitj-data-ingestion-buckets N. Virginia Shubham

iitj-data-ingestion-buckets info

Objects 7 Info

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

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	calendar_original.csv	csv	December 14, 2024, 15:51:06 (UTC+05:30)	171.7 MB	Standard
<input type="checkbox"/>	calendar.csv	csv	December 14, 2024, 15:51:06 (UTC+05:30)	140.7 MB	Standard
<input type="checkbox"/>	listings_prginal.csv	csv	December 14, 2024, 15:51:06 (UTC+05:30)	45.8 MB	Standard
<input type="checkbox"/>	listings.csv	csv	December 14, 2024, 15:51:08 (UTC+05:30)	1.7 MB	Standard
<input type="checkbox"/>	micro_merged.csv	csv	December 14, 2024, 15:51:10 (UTC+05:30)	5.0 MB	Standard
<input type="checkbox"/>	reviews_original.csv	csv	December 14, 2024, 15:51:06 (UTC+05:30)	93.2 MB	Standard
<input type="checkbox"/>	reviews.csv	csv	December 14, 2024, 15:51:06 (UTC+05:30)	93.6 MB	Standard

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

us-east-1.console.aws.amazon.com/s3/buckets?region=us-east-1&bucketType=general

AWS RDS

Amazon S3 > Buckets

Amazon S3

General purpose buckets

- Directory buckets
- Table buckets [New](#)
- Access Grants
- Access Points
- Object Lambda Access Points
- Multi-Region Access Points
- Batch Operations
- IAM Access Analyzer for S3

Block Public Access settings for this account

Storage Lens

- Dashboards
- Storage Lens groups
- AWS Organizations settings

Feature spotlight [10](#)

AWS Marketplace for S3

Account snapshot - updated every 24 hours [All AWS Regions](#)

Storage lens provides visibility into storage usage and activity trends. Metrics don't include directory buckets. [Learn more](#)

[View Storage Lens dashboard](#)

General purpose buckets [Info](#) [All AWS Regions](#)

Buckets are containers for data stored in S3.

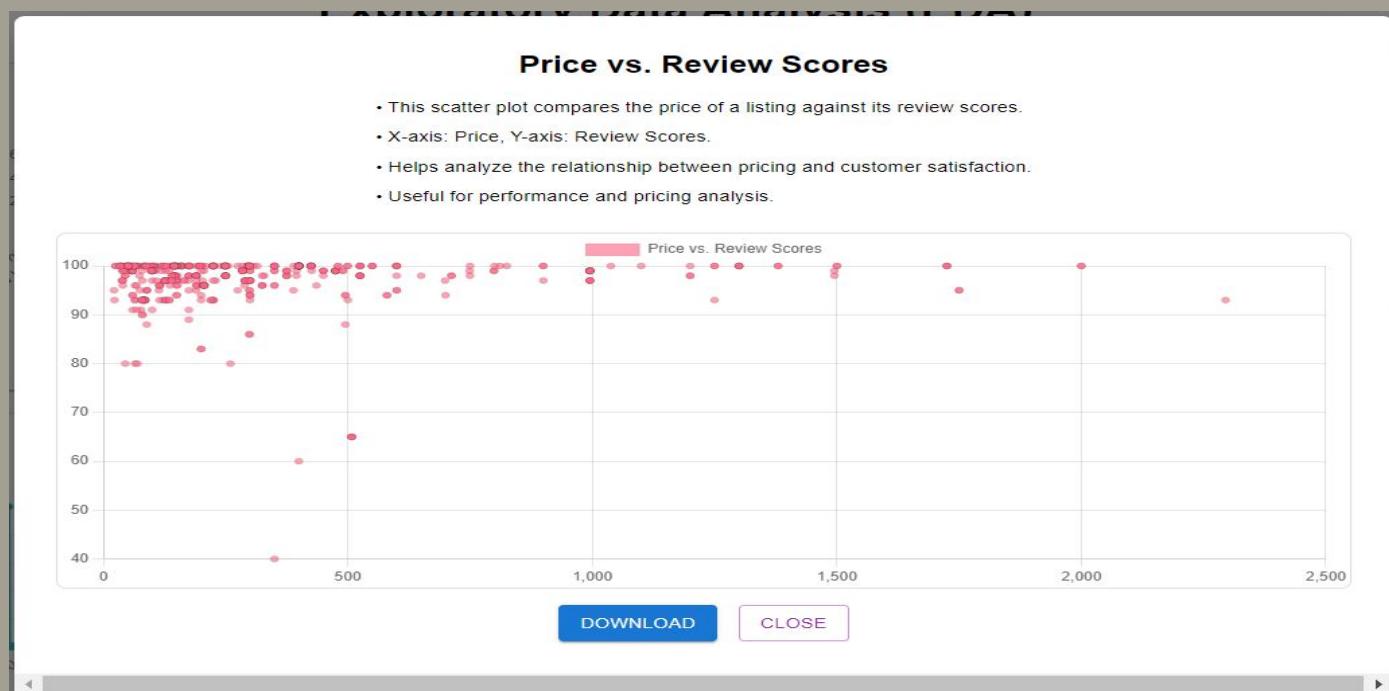
Find buckets by name

Name	AWS Region	IAM Access Analyzer	Creation date
iitj-data-ingestion-buckets	US East (N. Virginia) us-east-1	View analyzer for us-east-1	December 14, 2024, 14:29:33 (UTC+05:30)

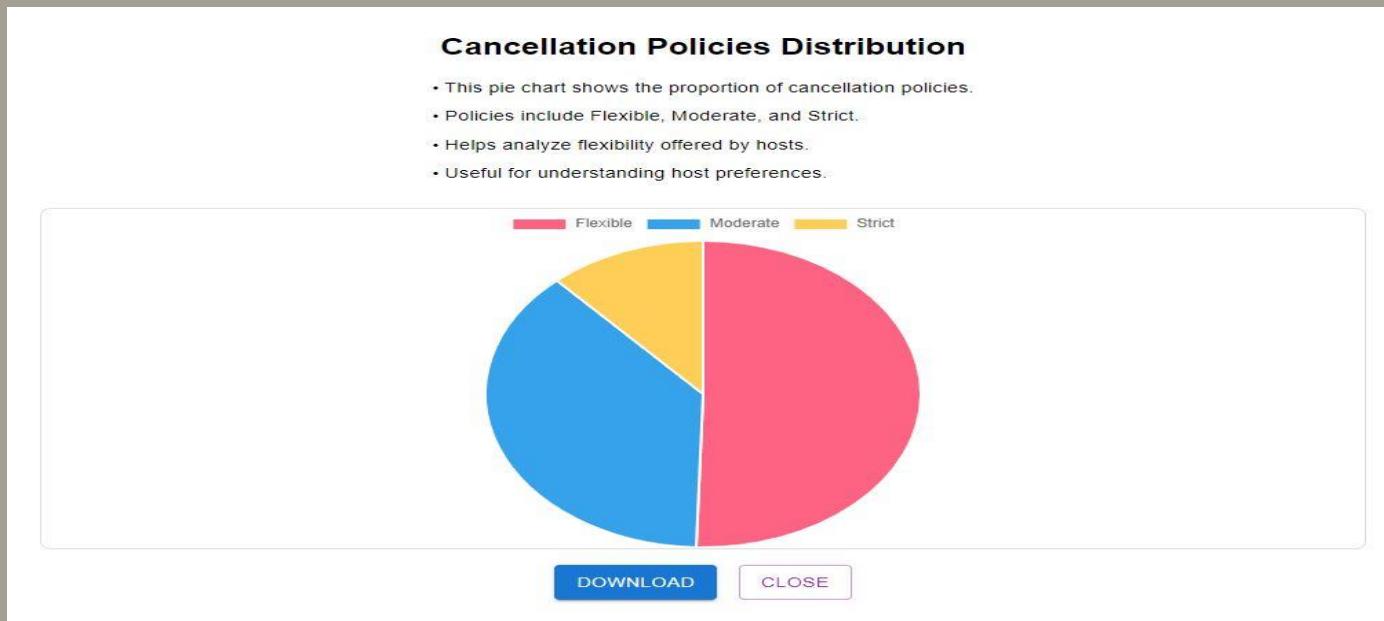
[Copy ARN](#) [Empty](#) [Delete](#) [Create bucket](#)

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

SCREENSHOTS FROM WEBSITE HOSTED



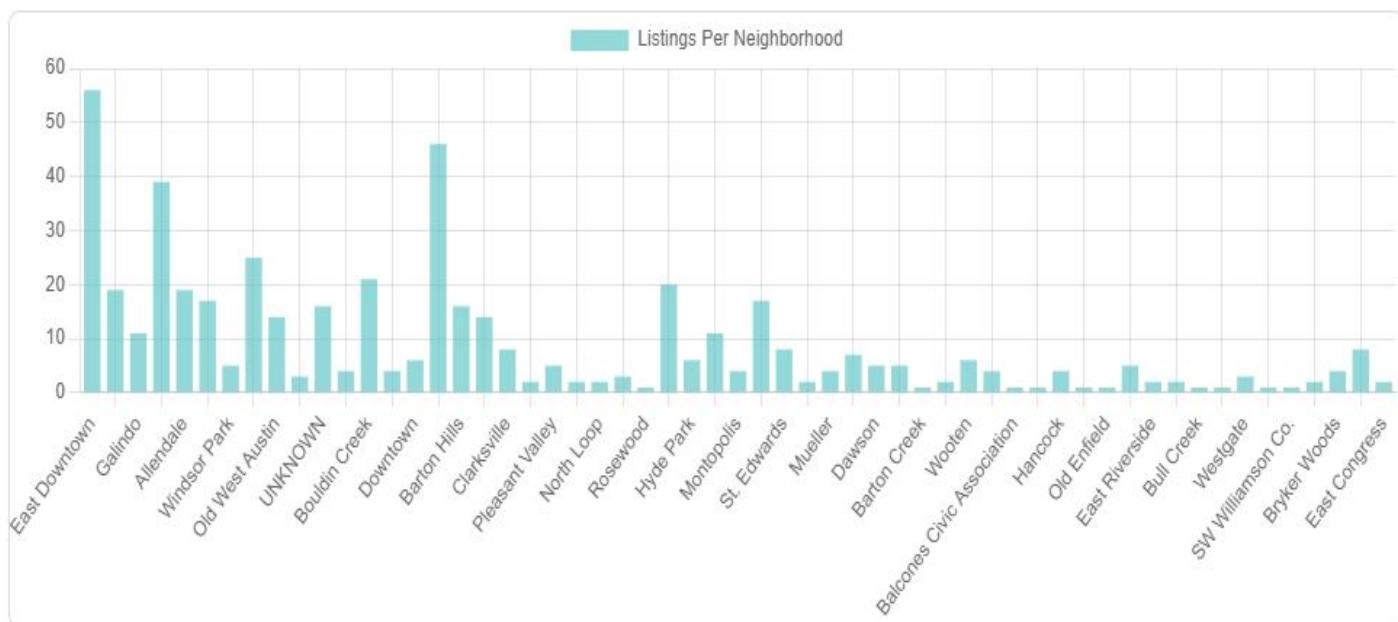
SCREENSHOTS FROM WEBSITE HOSTED



SCREENSHOTS FROM WEBSITE HOSTED

Listings Per Neighborhood

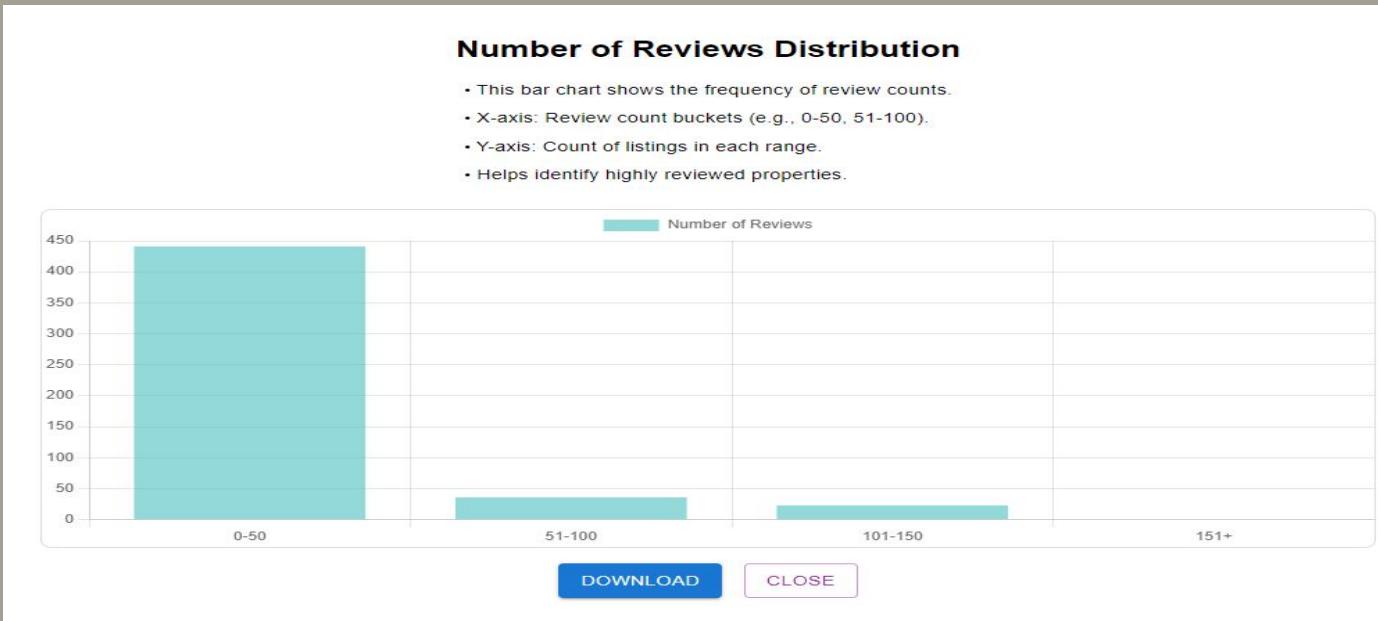
- This bar chart shows the count of listings per neighborhood.
- Helps identify the most popular neighborhoods.
- Neighborhood names are aggregated for clarity.
- Useful for analyzing the concentration of properties.



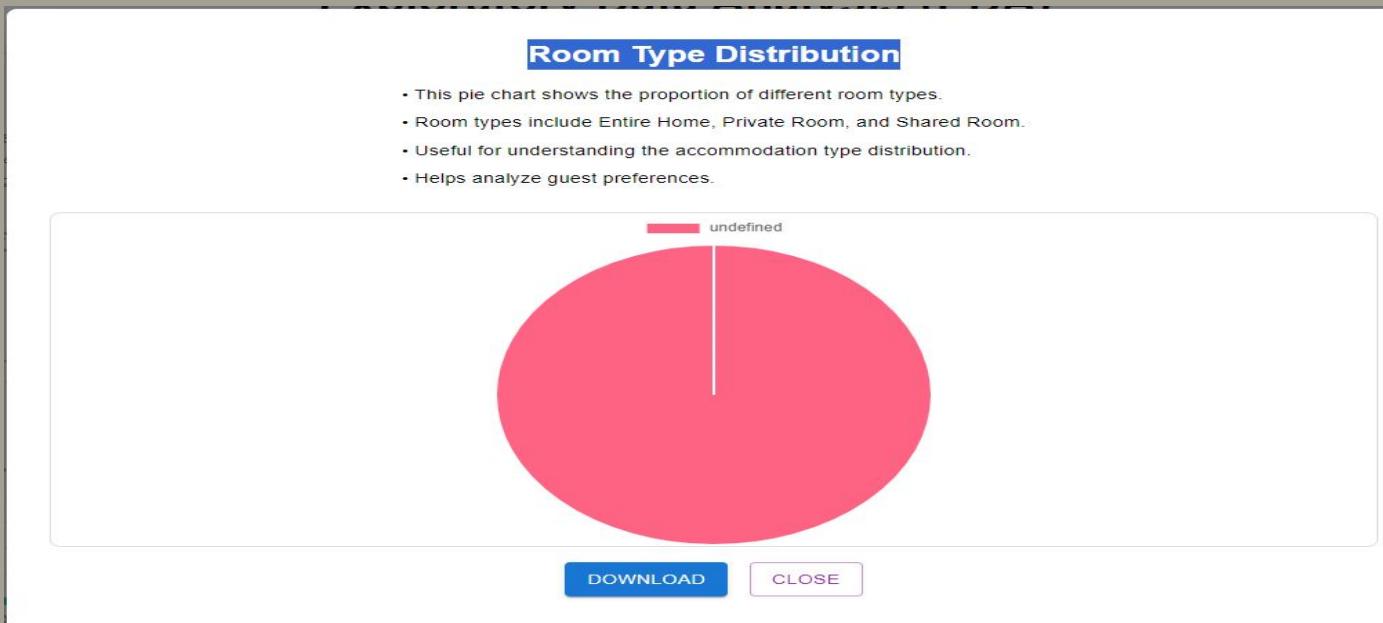
DOWNLOAD

CLOSE

SCREENSHOTS FROM WEBSITE HOSTED



SCREENSHOTS FROM WEBSITE HOSTED



SCREENSHOTS FROM WEBSITE HOSTED

IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis 1:49:57 PM

Unique Values Selection

AVAILABLE	LISTING URL
Select available	Select listing_url
CANCELLATION POLICY ID	LONGITUDE
Select cancellation_policy_id	Select longitude
COMMENTS	MAXIMUM NIGHTS
Select comments	Select maximum_nights
DATE	MINIMUM NIGHTS
Select date	Select minimum_nights
HOST ID	NAME
Select host_id	Select name
HOST NAME	NEIGHBOURHOOD

Contributors: Shubham Raj Bhagchandani Niraj Bhavesh Arora Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA

IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis 3:16:38 PM

Make Your Selection



Airbnb



Customer



Hotel Owner

Contributors: Shubham Raj Bhagchandani Niraj Bhavesh Arora Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA



FUTURE WORK RECOMMENDATION & CONCLUSION

FUTURE WORK AND RECOMMENDATION

Advanced Predictive Modeling

Future Work:

Develop sophisticated models (e.g., Random Forest, Neural Networks) to enhance predictions.

Recommendation:

Incorporate external data (local events, weather) for better accuracy in occupancy forecasting.

Dynamic Pricing Models

Future Work:

Create real-time pricing algorithms using market demand, competitor pricing, and events.

Recommendation:

Adopt dynamic pricing tools to maximize revenue and stay competitive in the market.

Sentiment Analysis of Reviews

Future Work:

Implement advanced NLP techniques to identify trends in guest sentiment and review themes.

Recommendation:

Improve guest satisfaction by addressing common concerns (cleanliness, communication).

Geospatial Analysis

Future Work:

Conduct geospatial analysis to uncover location-based trends and growth opportunities.

Recommendation:

Focus on high-demand neighborhoods for strategic property investments.

Longitudinal Studies

Future Work:

Analyze market trends over time and the impact of external factors like economic downturns.

Recommendation:

Use analytics tools to monitor key performance metrics and adjust strategies accordingly.

CONCLUSION

Dynamic Pricing Optimization:

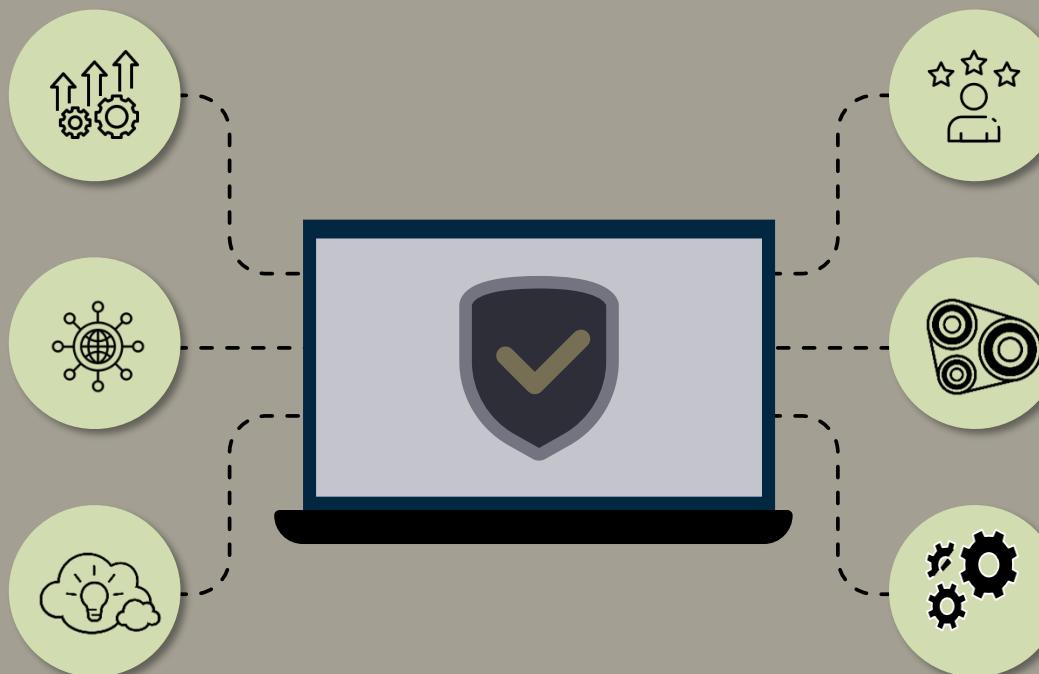
Data-driven pricing models have shown a significant correlation between pricing and occupancy rates. By analysing market demand and competitor pricing, hosts can optimize prices and improve occupancy by up to 15-20%.

Guest Sentiment Analysis:

Sentiment analysis of reviews reveals that listings with positive sentiment scores have 30% higher occupancy rates. Hosts can leverage NLP techniques to identify and address common guest concerns to improve ratings.

Room Type Analysis & Segmentation:

Data shows that entire homes have occupancy rates 25-30% higher compared to private or shared rooms, making them a preferred choice for families and larger groups.



Geospatial Data for Location Strategy:

Geospatial analysis indicates that certain neighbourhoods see up to 40% higher demand for rentals, highlighting the value of location-based insights for hosts looking to invest in emerging areas.

Predictive Analytics for Occupancy Forecasting:

Predictive models, utilizing features like price, review scores, and room type, can forecast occupancy rates with an accuracy of over 85%, enabling better planning and revenue predictions.

Continuous Data Monitoring & Real-Time Adjustments:

Hosts using real-time analytics dashboards experience a 20-25% improvement in performance metrics, allowing for rapid adjustments to pricing and availability, driving increased bookings.



Login

Login ID

Password

LOGIN



© 2024 IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

Contributors:

Shubham Raj

Bhagchandani Niraj Bhavesh Arora

Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA





Indian Institute of Technology
Jodhpur

Date: 12/17/2024

Project Report on

AIRBNB DATA CLEANING,
PROCESSING, AND ADVANCED
ANALYSIS

Welcome to AIRBNB DATA CLEANING, PROCESSING, AND ADVANCED ANALYSIS

START



© 2024 IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

Contributors:

Shubham Raj

Bhagchandani Niraj Bhavesh Arora

Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA





IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

3:32:19 PM

Individual Data Sets after performing (Step I & Step II)

[CALENDAR](#) [LISTINGS](#) [REVIEWS](#)

Ro...
5 ▾[SHOW MERGED DATA](#)

AVAILABLE	DATE	LISTING_ID	MAXIMUM_NIGHTS	MINIMUM_NIGHTS	PRICE
Yes	2019-09-19	320435	7	4	625
No	2019-09-19	320161	365	30	45
No	2019-09-20	320161	365	30	45
No	2019-09-21	320161	365	30	45
No	2019-09-22	320161	365	30	45

< 1 2 >

**Contributors:**

Shubham Raj



Bhagchandani Niraj Bhavesh Arora



Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA





IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

3:32:47 PM

Consolidated Dataset(Step III)

Record Count: 11136

Search

Row
5

LOAD DATA FROM RDS

LOAD DATA TO RDS

AVAILABLE	CANCELLATION POLICY ID	COMMENTS	DATE	HOST ID	HOST NAME	HOST SINCE	ID	ID X	ID Y	INSTANT BOOKABLE	LATITUDE	LISTING ID	LISTING URL	LONGITU
-----------	------------------------	----------	------	---------	-----------	------------	----	------	------	------------------	----------	------------	-------------	---------

0	0	The house is very nice. The metro/train station to downtown is a 10 min walk. When the weather is OK it's also a nice walk to downtown (40 minutes). A tip for the ..	2018-03-14	2466	Paddy	2008-08-23	22271	243130640	5245	0	30.27577	5245	https://www.airbnb.com/rooms/5245	-97.7137
---	---	---	------------	------	-------	------------	-------	-----------	------	---	----------	------	---	----------

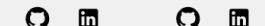


Contributors:

Shubham Raj



Bhagchandani Niraj Bhavesh Arora



Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA



Not secure 3.222.77.245:5173/hotel-owner-form

IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

3.222.77.245:5173 says
Record appended successfully!

ID X: deuy32f6glf

ID Y: fq9xj3vax18

INSTANT BOOKABLE: Instant Bookable

LATITUDE: 40.7128

LONGITUDE: -74.0060

NEIGHBOURHOOD: Downtown

NUMBER OF REVIEWS: 10

REVIEW SCORES RATING: 90

REVIEWER ID: reviewer123

REVIEWER NAME: Default Reviewer

ROOM TYPE ID: Entire Place

TYPE: Hotel Owner

SUBMIT

Contributors:

- Shubham Raj
- Bhagchandani Niraj
- Bhavesh Arora
- Jai Singh Kushwah
- JATIN SHRIVAS

PARAS PANDA

Elements Console Sources Network Performance

Fetch/XHR Doc CSS JS Font Img Media Manifest WS Wasm Other

Name: append-end

Request Payload

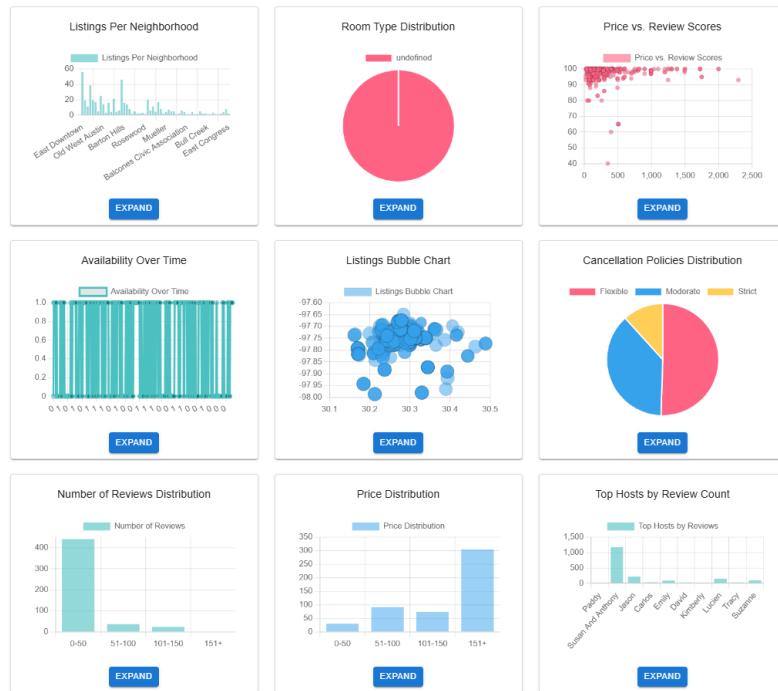
```
{listing_id: "89898912", name: "IIT Jodhpur", available: "1", cancellation_policy_id: "0", comments: "Default comments go here.", date: "2024-12-17", host_id: "iitj123", host_name: "Default Host", host_since: "2020-01-01", id_x: "deuy32f6glf", id_y: "fq9xj3vax18", instant_bookable: "1", latitude: "40.7128", longitude: "-74.0060", maximum_nights: "30", minimum_nights: "1", name: "IIT Jodhpur", neighbourhood: "Downtown", number_of_reviews: "10", price: "2500", review_scores_rating: "90", reviewer_id: "reviewer123", reviewer_name: "Default Reviewer", room_type_id: "2"}
```

2 requests | 316 B transferred | 77 E

IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

3:33:56 PM

Visualisation from RDS

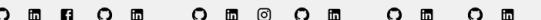


BACK

NEXT

Contributors:
Shubham Raj

Bhagchandani Niraj Bhavesh Arora Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA



Not secure

3.222.77.245:5173/fetch-unique-value



IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

Categorical Data Distribution

Select cancellation policy id

COMMENTS

Select comments

DATE

Select date

HOST ID

Select host id

HOST NAME

Select host name

HOST SINCE

Select host since

ID

Select id

ID X

Select id x

ID Y

Select id y

INSTANT BOOKABLE

Select instant bookable

LATITUDE

Select latitude

LISTING ID

Select longitude

MAXIMUM NIGHTS

Select maximum nights

MINIMUM NIGHTS

Select minimum nights

NAME

Select name

NEIGHBOURHOOD

Select neighbourhood

NUMBER OF REVIEWS

Select number of reviews

PRICE

Select price

REVIEW SCORES RATING

Select review scores rating

REVIEWER ID

Select reviewer id

REVIEWER NAME

Select reviewer name

ROOM TYPE ID

Select room type id

Contributors:

Shubham Raj

Bhagchandani Niraj Bhavesh Arora

Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA



IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

3:36:02 PM



You have selected Hotel Owner Form

LISTING ID
12345

NAME
Luxury Apartment

AVAILABLE
Available

PRICE
250

MINIMUM NIGHTS
1

MAXIMUM NIGHTS
30

CANCELLATION POLICY ID
No Policy

COMMENTS
Default comments go here.

DATE
17-12-2024

HOST ID
host123

HOST NAME
Default Host

HOST SINCE
2020-01-01

ID X
txkkg7k5e

ID Y
vh50ueyaf

INSTANT BOOKABLE
Not Instant Bookable

LATITUDE
40.7128

LONGITUDE
-74.0060

NEIGHBOURHOOD
Downtown

Contributors:
Shubham Raj Bhagchandani Niraj Bhavesh Arora Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA



IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

3:37:27 PM

Consolidated Dataset (After Adding New Data)

Record Count: 11136

Search

Row
5

LOAD DATA FROM RDS

LOAD DATA TO RDS

AVAILABLE	CANCELLATION POLICY ID	COMMENTS	DATE	HOST ID	HOST NAME	HOST SINCE	ID	ID X	ID Y	INSTANT BOOKABLE	LATITUDE	LISTING ID	LISTING URL	LONGITUDE	MAXIMUM NIGHTS	MINIMUM NIGHTS
0	2	We planned a last minute trip to Austin and came across Tina's place. We weren't really expecting much but a place to sleep since we would be out and about, but we loved Tina's	2018-06-14	193075619	Tina	2018-06-02	33386	276732119	25621580	0	30.28441	25621580	https://www.airbnb.com/rooms/25621580	-97.75965	365	1

< 1 ... 2223 2224 2225 ... 2228 >

BACK

NEXT

DOWNLOAD JSON



Contributors:

Shubham Raj

Bhagchandani Niraj Bhavesh Arora

Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA





IITJ - Airbnb Data Cleaning, Processing and Advanced Analysis

3:38:49 PM

Consolidated Dataset (After Adding New Data)

Record Count: 11136

Search –
niraj

Row
5

LOAD DATA FROM RDS

LOAD DATA TO RDS

AVAILABLE	CANCELLATION POLICY ID	COMMENTS	DATE	HOST ID	HOST NAME	HOST SINCE	ID	ID X	ID Y	INSTANT BOOKABLE	LATITUDE	LISTING ID	LISTING URL	LONGITUDE	MAXIMUM NIGHTS	MINIMUM NIGHTS	NAME	NEIGHBOURHOOD
1	0	Default comments go here.	2024-12-17	host123	Default Host	2020-01-01	33404	chec4z3pwzt	en27g5urude	0	40.7128	999999999999	-	-74.0060	30	1	Niraj Apartment	Downtown

< 1 >

BACK

DOWNLOAD JSON

NEXT



Contributors:

Shubham Raj



Bhagchandani Niraj JAI SINGH KUSHWAH





DATA COOKING IN PROGRESS...



BACK

NEXT



Contributors:

Shubham Raj

Bhagchandani Niraj Bhavesh Arora

Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA





Future Room Price Prediction

NAME *	YEAR *
<input type="text" value="The Virgo Den"/>	<input type="text" value="2025"/>
DAY OF WEEK *	MINIMUM NIGHTS
<input type="text" value="Monday"/>	<input type="text" value="1"/>
MAXIMUM NIGHTS	NUMBER OF REVIEWS
<input type="text" value="365"/>	<input type="text" value="1"/>
REVIEW SCORES RATING	ROOM TYPE ID
<input type="text" value="100"/>	<input type="text" value="0"/>
CANCELLATION POLICY ID	AVAILABLE
<input type="text" value="2"/>	<input type="text" value="0"/>
INSTANT BOOKABLE	
<input type="text" value="0"/>	
<input type="button" value="SUBMIT"/>	

BACK

NEXT

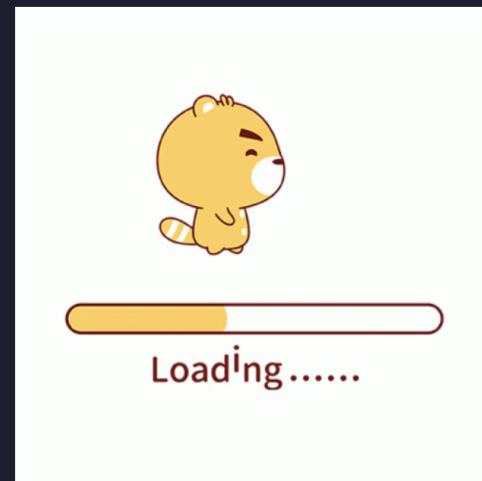
**Contributors:**

Shubham Raj Bhagchandani Niraj Bhavesh Arora Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA





- Filling in the data...**
Preparing input for the model
- Training the data...**
Building and refining the model
- Predicting the data...**
Generating the final output





Hotel "Private Master Bedroom & Bath Oasis in Lovely Home" on "Thursday" for year "2028" is

\$94.72

name: Private Master Bedroom & Bath Oasis in Lovely Home	minimum nights: 9	maximum nights: 15
number of reviews: 4	review scores rating: 100	room type id: 1
cancellation policy id: 1	year: 2028	day of week: Thursday
available: 0	instant bookable: 1	comments: The host canceled this reservation 23 days before arrival. This is an automated posting.
date: 2018-06-11	host id: 49374084	host name: Kendall
host since: 2015-11-18	id: 33401	id x: 275657632
id y: 25842890	latitude: 30.27235	listing id: 25842890
listing url: https://www.airbnb.com/rooms/25842890	longitude: -97.71527	neighbourhood: Rosewood
price: 108	reviewer id: 194201346	reviewer name: Norma

Re-predictNext BI Graphs



Page Not Found

Redirecting to the Home Page in 10 seconds...

[Go to Home](#)



Contributors:

Shubham Raj

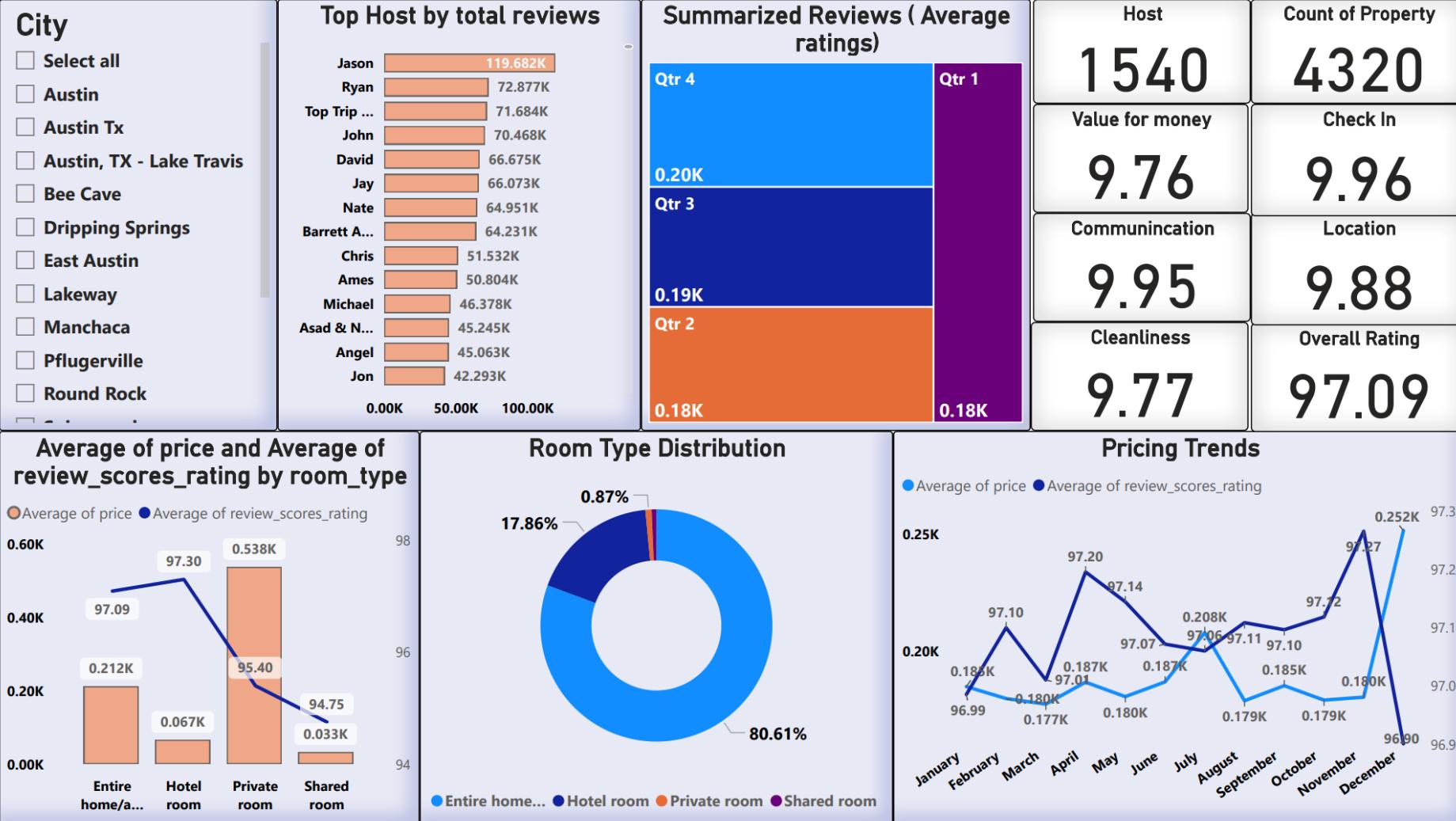


Bhagchandani Niraj Bhavesh Arora



Jai Singh Kushwah JATIN SHRIVAS PARAS PANDA







THANK YOU !!!

