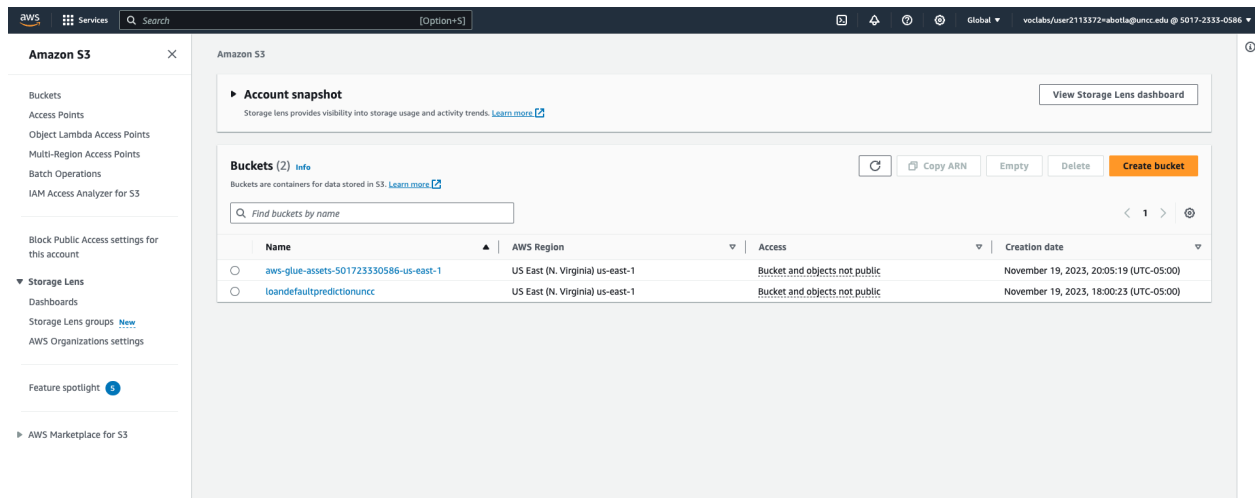# Project Deliverable 2
## Team - 4
## Abhinav Botla
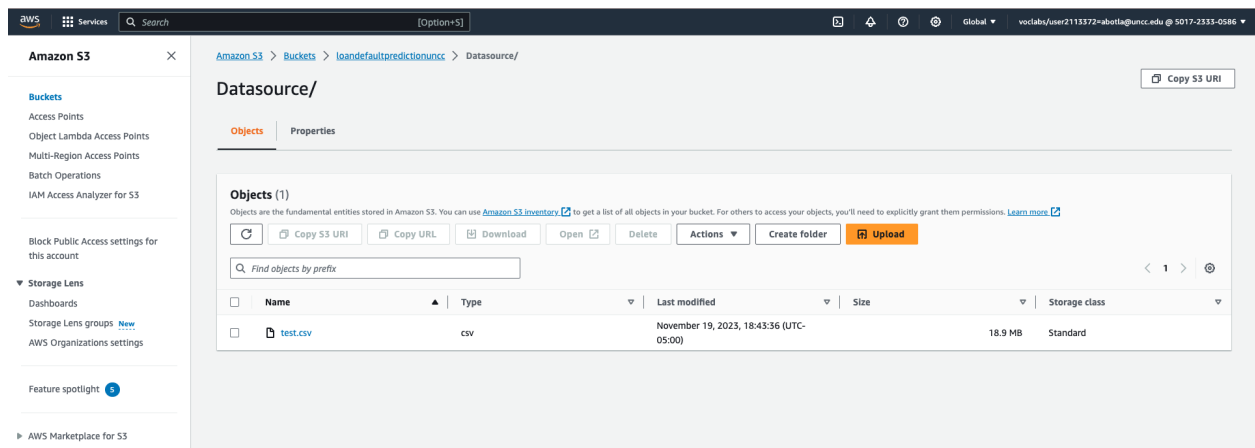## Ajith Gannamaneni
## Hrushikesh Dandge
## Mitra Buggaveeti
## Prateek Chanddra

Created the S3 bucket with name of **loandefaultpredictionuncc** to store the project data in the bucket



Created the Datasource folder in the bucket to store the data

Analysis of data by using the AWS Athena:

Creating table DDL in Athena with predefined attributes already available in the data source.



Time based exploratory data analysis:

Analyzing the credit history of the customers:



Analyzing the count of customers with specific CNS Score:

Athena now supports typeahead code suggestions to speed up SQL query development
Typeahead suggestions are turned on by default. You can change this setting in query editor preferences.

Edit preferences

**Data**

Query 3 | Query 4 | Query 5 | Query 6 | Query 7 | Query 8 | Query 9

Data source
AwsDataCatalog

Database
loan_data

```
1  SELECT
2    AVG(disbursed_amount) AS avg_disbursed_amount,
3    MIN(disbursed_amount) AS min_disbursed_amount,
4    MAX(disbursed_amount) AS max_disbursed_amount
5  FROM loan_data.test;
6
```

Tables and views | Create

Filter tables and views

Tables (1)                1
  test

Views (0)                 1

SQL    Ln 1, Col 1

Run again | Explain | Cancel | Clear | Create

Reuse query results
up to 60 minutes ago

Query results | Query stats

✓ Completed          Time in queue: 97 ms    Run time: 685 ms    Data scanned: 18.95 MB

**Results (1)**          Copy | Download results

Search rows                                                     1

| # | avg_disbursed_amount | min_disbursed_amount | max_disbursed_amount |
|---|---|---|---|
| 1 | 56076.80326891594 | 11613 | 940690 |

**AWS GLUE:**
**CRAWLERS:**

**AWS Glue**

Getting started
ETL jobs
  Visual ETL
  Notebooks
  Job run monitoring
Data Catalog tables
Data connections
Workflows (orchestration)
Data Catalog
  Databases
  Tables
Stream schema registries
  Schemas
Connections
Crawlers
  Classifiers
Catalog settings
Data Integration and ETL
  ETL jobs
  Visual ETL
  Notebooks
  Job run monitoring
Interactive Sessions
Data classification tools
  Sensitive data detection
  Record Matching
Triggers
Workflows (orchestration)
  Blueprints
Security configurations
Legacy pages

AWS Glue > Crawlers

**Crawlers**

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

**Crawlers (1)** Info          Last updated (UTC)
View and manage all available crawlers.   November 20, 2023 at 01:41:21

Filter crawlers          Action | Run | **Create crawler**

| | Name | State | Schedule | Last run | Last run timestamp | Log | Table changes from last r... |
|---|---|---|---|---|---|---|---|
| | gluecrawler | ✓ Ready | | ✓ Succeeded | November 20, 2023 at 24:... | View log | 1 updated |

Glue ETL Pipeline:

## Visualization of the data flow



We doesnt pick the records which have a disbursed amount below 15000.

Dropping the duplicates from the data



Created the new column in the data set with attribute name of Available amount , which gives the information about (sanctioned amount - disbursed amount)

Final destination:



Not able to access the Quick Sight.
Quick sight