# Project Report on ETL Pipeline (ADF)

**Project Overview**

The project aims to create an efficient data engineering pipeline to process loan data from GitHub, transform it for analytics, and visualize insights using Power BI. The pipeline involves Azure Data Factory (ADF) for ETL (Extract, Transform, Load) operations, Azure Synapse Analytics for staging and transformation, and Power BI for reporting and visualization.

**Tools and Technologies**

- **Data Source:** GitHub (Loan Dataset)
- **ETL Platform:** Azure Data Factory (ADF)
- **Data Storage:** Azure Blob Storage and Azure Synapse Analytics
- **Visualization Tool:** Power BI
- **Data Transformation:** Data Flow activity in ADF
- **Validation:** Script activity in ADF using SQL scripts

**Key Layers**

1. **Raw Data Layer:** Raw data extracted from GitHub and stored in Azure Blob Storage.
2. **Staging Layer:** Data loaded into Azure Synapse Analytics for validation and temporary storage.
3. **Processed Data Layer:** Transformed and validated data stored in Synapse's final table for analytics.
4. **Presentation Layer:** Data fetched into Power BI for generating dashboards and reports.

**Data and File Formats**

- **Source Format:** CSV files from GitHub
- **Intermediate Format:** Parquet (for transformation in ADF)
- **Final Table Schema:** Synapse table LoanProcessData with detailed loan attributes.

**Storage Location**

- **Raw Data:** Azure Blob Storage
- **Transformed Data:** Azure Synapse Analytics (staging and processed tables)
- **Visualization Data:** Synapse Analytics integrated with Power BI

**Dataset Understanding:**

**1. Unique Identifiers**

- **id**: Critical.

  **Reason**: Identifies each loan application. If missing, the row loses uniqueness, making it impossible to track specific loans.

  **Action**: Check for null or duplicate values.

- **member_id**: Critical.

  **Reason**: Identifies the borrower. Without this, customer-level analysis cannot be performed.

  **Action**: Check for null values.

---

## 2. Loan Attributes

- **loan_amount**: Critical.

  **Reason**: Core financial figure. Null values affect financial metrics such as total disbursed loans.

  **Action**: Check for null or zero values.

- **int_rate (Interest Rate)**: Important.

  **Reason**: Used in calculating the cost of borrowing. Null values would mislead analysis on loan profitability.

- **installment**: Important.

  **Reason**: Monthly repayment amount. Null values impact cash flow predictions.

- **term**: Important.

  **Reason**: Represents loan duration (e.g., 36 months). Null values would prevent understanding loan periods.

---

## 3. Borrower Information

- **application_type**: Important.

  **Reason**: Indicates individual vs. joint application. Helps in customer segmentation.

- **emp_length (Employment Length)**: Important.

  **Reason**: Used to evaluate borrower stability. Null or missing values reduce reliability of risk assessments.

- **emp_title (Employment Title)**: Less important (depends on analysis goal).

  **Reason**: Provides job information. Not critical for financial calculations but useful in customer profiling.

- **annual_income**: Critical.

  **Reason**: Key determinant for loan approval and debt-to-income (DTI) ratio calculations.

- **dti (Debt-to-Income)**: Important.

  **Reason**: Measures borrower's debt relative to income. Null values affect risk evaluation.

---

## 4. Loan Status and Dates

- **loan_status**: Critical.

  **Reason**: Determines whether the loan is current, defaulted, or paid off. Null values prevent tracking performance.

- **issue_date**: Critical.

  **Reason**: Date loan was issued. Needed for time-based analysis.

- **last_payment_date**: Important.

  **Reason**: Helps analyze payment history. Missing values are acceptable if the loan is new.

- **next_payment_date**: Important.

  **Reason**: Relevant for forecasting future payments. Null values are acceptable if the loan is paid off.

- **last_credit_pull_date**: Less critical.

  **Reason**: Date of last credit check.

---

**5. Credit Information**

- **total_acc (Total Accounts)**: Important.

  **Reason**: Number of credit accounts affects creditworthiness.

- **total_payment**: Critical.

  **Reason**: Represents the total payments made. Affects revenue and performance calculations.

---

**6. Risk Assessment**

- **grade** and **sub_grade**: Important.

  **Reason**: Reflect credit risk levels.

- **Good vs Bad Loan**: Critical.

  **Reason**: Directly indicates loan performance. Null values undermine model training for risk prediction.

---

**7. Other Attributes**

- **home_ownership**: Less critical.

  **Reason**: Provides collateral information.

- **verification_status**: Important.

  **Reason**: Indicates if income was verified.

- **purpose**: Important.

  **Reason**: Describes the loan purpose. Useful for segmentation.

- **address_state**: Less critical.

  **Reason**: Geographic information for loan distribution.

**Architectural Diagram:**

**Data Engineering Pipeline Design:**



1. **Copy Data Activity:**
    - ○ **Objective:** Extract raw data from GitHub and load it into Azure Blob Storage.
    - ○ **Configuration:**
        - ▪ Source: GitHub dataset.
        - ▪ Sink: Azure Blob Storage (Raw Data Container).

**Storage of Data to bronze:**



2. **Data Flow Activity:**



- o **Objective:** Transform data and load it into the Synapse final table.

- o **Transformation Operations:**

    - ▪ Data cleaning (null handling and type conversion).

### 3. Summary of Critical Columns to Check for Null/Empty:

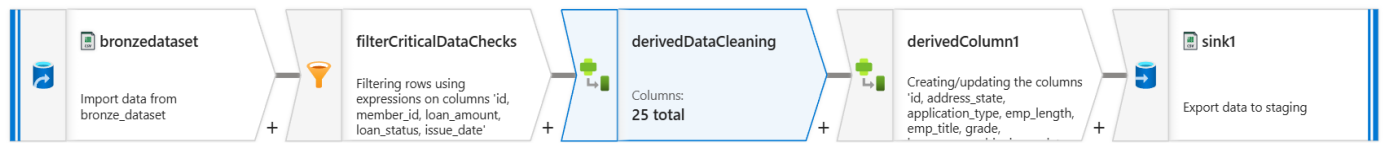| Column Name | Why It's Critical |
|---|---|
| id | Unique identifier for loans |
| member_id | Unique identifier for borrowers |
| loan_amount | Core financial metric |
| annual_income | Determines borrower's financial capability |
| loan_status | Tracks loan performance |
| issue_date | Needed for time-based analysis |
| Good vs Bad Loan | Indicates loan risk and performance |
| total_payment | Represents total payments made |



## Columns Important for Data Quality:

| Column Name | Why It's Important |
|---|---|
| emp_length | Employment stability |
| dti | Debt-to-income ratio |
| int_rate | Loan cost |
| term | Loan duration |
| grade, sub_grade | Credit risk assessment |

| loan_amount | ∨ | toFloat(loan_amount) | 1.2f | + | 🗑 |
| loan_status | ∨ | upper(loan_status) | abc | + | 🗑 |
| total_payment | ∨ | toFloat(total_payment) | 1.2f | + | 🗑 |
| issue_date | ∨ | toDate(issue_date, 'dd-MMM-yy') | 📅 | + | 🗑 |

- o **Output:** Transformed data written to LoanProcessData table.

4. **Copy Data Activity:**

   - o **Objective:** Load clean data from Blob Storage to Synapse staging table.

   - o **Configuration:**

     - ▪ Source: Azure Blob Storage.

     - ▪ Sink: Synapse staging table.

**Storage of clean silver Data:**



5. **Script Activity:**

   o **Objective:** Validate transformed data.

   o **Validation Checks:**

      ▪ Duplicate records.

      ▪ Null values in critical columns.

**Synapse Data Quality Check:**



6. **Power BI Integration:**

   o **Objective:** Visualize data from Synapse final table.

   o **Configuration:**

      ▪ Data fetched using DirectQuery for real-time insights.

      ▪ Dashboards include metrics like loan status distribution, average interest rates, and bad loan percentages.

**Integration with Power BI**

Power BI connects directly to the Synapse table "LoanOProcessData" for real-time reporting. The dashboards provide actionable insights such as:

- Loan status analysis ("Good Vs Bad Loan" split).

- Top states by loan amount.

- Average interest rates by grade and sub-grade.

## DataModelling:

**Dashboard:**



## BANK LOAN REPORT | SUMMARY

| | | | | |
|---|---|---|---|---|
| **Total Loan Applications** 37K | **Total Funded Amount** $418.9M | **Total Received Amount** $455.6M | **Average Interest Rate** 12.1% | **Average DTI** 13.4% |
| MTD 4.1K / MoM 6.1% | MTD $51.0M / MoM 12.0% | MTD $55.1M / MoM 15.1% | MTD 12.4% / MoM 3.6% | MTD 13.8% / MoM 2.6% |

**GOOD LOAN ISSUED** — 86.4%

Good_Loan_Application: 32K
Good_Loan_Funded_Amount: $356.7M
Good_Loan_Recevied_Amount: $420.0M

**BAD LOAN ISSUED** — 13.6%

Bad_Loan_Application: 5K
Bad_Loan_Funded_Amount: $62.2M
Bad_Loan_Recevied_Amount: $35.5M

### LOAN STATUS

| loan_status | TotalLoanApplications | TotalFundedAmount | Total Amount Received | MTD Funded Amount | MTDTotalAmountReceived | Avg Interest Rate | Avg DTI |
|---|---|---|---|---|---|---|---|
| Charged Off | 5067 | $6,22,48,700 | $3,55,41,160 | $81,76,700 | $50,36,723 | 13.88% | 14.11% |
| Current | 1066 | $1,82,29,350 | $2,33,99,779 | $38,57,525 | $48,28,582 | 15.07% | 14.78% |
| Fully Paid | 31005 | $33,84,23,600 | $39,66,42,813 | $3,89,74,575 | $4,52,30,711 | 11.65% | 13.23% |
| Grand Total | 37138 | $41,89,01,650 | $45,55,83,752 | $5,10,08,800 | $5,50,96,016 | 12.05% | 13.40% |

**STATE**: All
**GRADE**: All
**PURPOSE**: All

Summary | Overview | Details

**Conclusion**

The project successfully implemented a robust ETL pipeline to process loan data from GitHub and provide valuable business insights through Power BI dashboards. The integration of ADF and Synapse Analytics ensures scalability and real-time data processing capabilities.

**Future Enhancements**

1. **Automated Scheduling:** Implementing triggers in ADF to automate the data pipeline execution.

2. **Real-Time Data Updates:** Integrating Event Hub for streaming real-time data changes.

3. **Advanced Analytics:** Using Synapse ML for predictive modeling of loan defaults.

4. **Data Governance:** Implementing data masking and role-based access control (RBAC) for sensitive information.

5. **Dashboard Enhancements:** Adding drill-through and cross-filtering capabilities in Power BI dashboards for deeper insights.