

10 Academy: Artificial Intelligence Mastery

Week 2 Challenge Document

Customer Experience Analytics for
Fintech Apps

A Real-World Data Engineering Challenge: Scraping, Analyzing, and Visualizing Google
Play Store Reviews

Date: 26 Nov - 02 Dec 2025

Challenge Overview

This week's challenge centers on analyzing customer satisfaction with mobile banking apps by collecting and processing user reviews from the Google Play Store for three Ethiopian banks:

- Commercial Bank of Ethiopia (CBE)
- Bank of Abyssinia (BOA)
- Dashen Bank

You'll scrape app reviews, analyze sentiments and themes, and visualize insights as a Data Analyst at Omega Consultancy, a firm advising banks.

Building on Week 1's foundational skills, this week introduces web scraping, thematic NLP analysis, and basic database engineering.

Business Objective

Omega Consultancy is supporting banks to improve their mobile apps to enhance customer retention and satisfaction. Your role as a Data Analyst is to:

- Scrape user reviews from the Google Play Store.
- Analyze sentiment (positive/negative/neutral) and extract themes (e.g., "bugs", "UI").
- Identify satisfaction drivers (e.g., speed) and pain points (e.g., crashes).
- Store cleaned review data in a Postgres database.
- Deliver a report with visualizations and actionable recommendations.

Scenarios

These scenarios simulate real consulting tasks faced by product, marketing, and engineering teams in the banking industry.

1. Scenario 1: Retaining Users

CBE has a 4.2, BOA 3.4, and Dashen 4.1 star rating. Users complain about slow loading during transfers. Analyze if this is a broader issue and suggest areas for app investigation.

2. Scenario 2: Enhancing Features

Extract desired features (e.g., transfer, fingerprint login, faster loading times) through keyword and theme extraction. Recommend how each bank can stay competitive.

3. Scenario 3: Managing Complaints

Cluster and track complaints (e.g., "login error") to guide AI chatbot integration and faster support resolution strategies.

These scenarios highlight user retention, feature innovation, and support efficiency, aligning with fintech priorities.

Dataset Overview

You will scrape reviews from the Google Play Store for three banks' apps, collecting:

- Review Text: User feedback (e.g., "Love the UI, but it crashes often").
- Rating: 1–5 stars.
- Date: Posting date (time optional).
- Bank/App Name: E.g., "Commercial bank of Ethiopia Mobile".
- Source: Google Play.

Minimum: 400 reviews per bank (1200 total). If scraping is limited,

Team

Facilitators:

- Kerod
- Mahbubah
- Filimon

Key Dates

- **Introduction:** 10:30 AM UTC, Wednesday, 26 Nov 2025
- **Interim Submission:** 8:00 PM UTC, Sunday, 30 Nov 2025
- **Final Submission:** 8:00 PM UTC, Tuesday, 02 Dec 2025

Communication & Support

- **Slack channel:** #all-week-2
- **Office hours:** Mon–Fri, 08:00–15:00 UTC

Learning Objectives

By the end of this week, you will be able to:

- Scrape and preprocess user reviews from the Google Play Store to prepare text data for analysis.
- Apply NLP techniques to analyze review sentiment and identify key customer feedback themes.
- Design and implement a relational database schema in Postgres to store and manage review data.
- Derive actionable insights from analyzed data and create compelling visualizations for business stakeholders.
- Develop and present a data-driven report with recommendations for app improvement in a fintech context.
- Employ Git for version control and write unit tests to ensure the reliability of data processing scripts.

Project Planning - Data Collection & Analysis

Tasks:

- Scrape and preprocess reviews.
- Analyze sentiment and themes.
- Synthesize insights and visualize results.

KPIs:

- **Proactivity:** Sharing scraping/NLP references.
- **Data Quality:** 1,200+ clean reviews with <5% errors.
- **Insights:** 3+ drivers/pain points per bank.
- **Clarity:** Stakeholder-friendly visualizations.

Deliverables and Tasks

Task 1: Data Collection and Preprocessing

Scrape reviews from the Google Play Store, preprocess them for analysis, and manage code via GitHub.

Tasks:

- Git Setup:
 - Create a GitHub repository.
 - Include .gitignore, requirements.txt.
 - Use "task-1" branch, Commit frequently with meaningful messages, ideally after completing logical chunks of work or at the end of each work session.
- Web Scraping:
 - Use google-play-scraper to collect reviews, ratings, dates, and app names for three banks.
 - Target a minimum of 400+ reviews per bank (1,200 total).
- Preprocessing:
 - Remove duplicates, handle missing data.
 - Normalize dates (e.g., to YYYY-MM-DD).
 - Save as CSV with columns: review, rating, date, bank, source.

KPIs:

- 1,200+ reviews collected with <5% missing data.
- Clean CSV dataset.
- Organized Git repo with clear commits.

Minimum Essential:

- Scrape at least 400 reviews per bank (1200 total).
- Commit a preprocessing script.
- Update README.md with methodology.

Task 2: Sentiment and Thematic Analysis

Description: Quantify review sentiment and identify themes to uncover satisfaction drivers and pain points.

Tasks:

- Sentiment Analysis:
 - Use [distilbert-base-uncased-finetuned-sst-2-english](#) to compute sentiment scores (positive, negative, neutral). Alternatively, you can start with simpler libraries like [VADER](#) or [TextBlob](#) and compare results if time permits.
 - Aggregate by bank and rating (e.g., mean sentiment for 1-star reviews).
- Thematic Analysis:
 - A theme refers to a recurring concept or topic within user reviews. For this challenge, themes will help summarize user feedback into actionable categories for the banks.
 - **Keyword Extraction & Manual/Rule-Based Clustering:**
 - Extract significant keywords and n-grams using TF-IDF or spaCy (e.g., "login error", "slow transfer", "good UI").
 - To aid in grouping these keywords and understanding broader review topics, you can optionally employ topic modeling techniques
 - Group related keywords and phrases into 3-5 overarching themes per bank (e.g., 'Account Access Issues', 'Transaction Performance', 'User Interface & Experience', 'Customer Support', 'Feature Requests'). Document your grouping logic.
 - Pipeline:
 - Script preprocessing (tokenization, stop-word removal, lemmatization if useful) with Pandas and NLP libraries.
 - Save results as CSV (e.g., review_id, review_text, sentiment_label, sentiment_score, identified_theme(s)).
 - Extract keywords with spaCy or TF-IDF (e.g., "crash", "support").
 - Cluster into 3-5 themes per bank (e.g., UI, reliability).
- Git:
 - Use "task-2" branch, commit scripts, merge via pull request.

KPIs:

- Sentiment scores for 90%+ reviews.
- 3+ themes per bank with examples.
- Modular pipeline code.

Minimum Essential:

- Sentiment scores for 400 reviews.
- 2 themes per bank via keywords.
- Commit analysis script.

Task 3: Store Cleaned Data in PostgreSQL

Design and implement a relational database in PostgreSQL to persistently store the cleaned and processed review data. This step simulates real-world data engineering

workflows, where PostgreSQL is widely used for its robustness, ease of use, and strong community support.

Tasks:

- PostgreSQL Database Setup
 - Install PostgreSQL on your system.
 - PostgreSQL is generally straightforward to install across different operating systems. If you encounter issues, consult facilitators or refer to the official documentation.
- Create a database named `bank_reviews`
- Define schema:
 - Banks Table: Stores information about the banks.
 - Suggested columns: bank_id (PRIMARY KEY), bank_name, app_name
 - Reviews Table: Stores the scraped and processed review data.
 - Suggested columns: review_id (PRIMARY KEY), bank_id (FOREIGN KEY), review_text, rating, review_date, sentiment_label, sentiment_score, source
- Insert cleaned review data using Python (psycopg2 or SQLAlchemy recommended)
- Write SQL queries to verify data integrity (e.g., count reviews per bank, average rating)

KPIs:

- Working database connection + insert script
- Tables populated with >1,000 review entries
- SQL dump or schema file committed to GitHub

Minimum Essential:

- PostgreSQL database created with both tables
- Python script that successfully inserts at least 400 reviews
- Schema documented in README.md

Task 4: Insights and Recommendations

Description: Derive insights from sentiment and themes, visualize results, and recommend app improvements.

Tasks:

- Insights:
 - Identify 2+ drivers (e.g., fast navigation) and pain points (e.g., crashes) per bank..
 - Compare banks (e.g., CBE vs. BOA).

- Suggest 2+ improvements (e.g., add budgeting tool) per bank..
- Visualization:
 - Create 3–5 plots (Matplotlib, Seaborn): sentiment trends, rating distributions, keyword clouds.
- Ethics:
 - Note potential review biases (e.g., negative skew).
- Git:
 - Use “task-4” branch, commit visuals/reports, merge via pull request.

KPIs:

- 2+ drivers/pain points with evidence.
- Clear, labeled visualizations.
- Practical recommendations.

Minimum Essential:

- 1 driver, 1 pain point per bank.
- 2 plots (e.g., sentiment bar, keyword chart).
- 4-page final report.

Due Dates

- **Sunday, 30 Nov 2025, 8:00 PM UTC:**
 - **GitHub Link:** Main branch with "task-1" merged, partial "task-2".
 - **Interim Report:** Maximum of 4 pages, covering scraping and early analysis.
- **Tuesday, 02 Dec 2025, 8:00 PM UTC:**
 - **GitHub Link:** Main branch with all tasks.
 - **Final Report:** Maximum of 10 pages and a maximum of 10 plots, Medium style. PDF format.

Other Considerations

- **Documentation:** Comment code, update README.md.
- **Collaboration:** Use GitHub Issues for coordination.
- **Professionalism:** Meet deadlines, learn proactively.
- **Flexibility:** Pivot to datasets if scraping fails.

Tutorials Schedule

In the following, the color **purple** indicates morning sessions, and non-purple indicates afternoon sessions.

- **Day 1:**
 - Challenge Intro & UX in Fintech (Kerod)
 - Scraping Basics and Data Cleaning and Preprocessing Strategy (Filimon)
- **Day 2:**
 - Sentiment Analysis and Keyword Extraction (Filimon)
 - Database Fundamentals & PostgreSQL Setup (Mahbubah)
- **Day 3:**
 - Error Handling (Mahbubah)
 - Introduction to Unit Testing (Semegn)
- **Day 4:**
 - Insight Communication and Plotting for Stakeholders (Mahbubah)
 - Q&A (Mahbubah)
- **Day 5:**
 - Q&A (Filimon)

References

- **Web Scraping (Google Play Store):**
 - google-play-scraper library:
<https://pypi.org/project/google-play-scraper/>
 - GitHub Repository: <https://github.com/JoMingyu/google-play-scraper>
- **Data Handling & Preprocessing:**
 - Pandas Documentation:
 - 10 Minutes to pandas:
https://pandas.pydata.org/pandas-docs/stable/user_guide/10min.html
 - Working with Text Data:
https://pandas.pydata.org/pandas-docs/stable/user_guide/text.html
 - Handling Missing Data:
https://pandas.pydata.org/pandas-docs/stable/user_guide/missing_data.html
- **Sentiment Analysis & NLP:**
 - Hugging Face transformers (for distilbert-base-uncased-finetuned-sst-2-english):
 - Documentation - Pipelines (Sentiment Analysis):
https://huggingface.co/docs/transformers/main_classes/pipelines#transformers.TextClassificationPipeline
 - Model Card (distilbert-base-uncased-finetuned-sst-2-english):
<https://huggingface.co/distilbert-base-uncased-finetuned-sst-2-english>
 - VADER Sentiment :
 - GitHub Repository: <https://github.com/cjhutto/vaderSentiment>
 - TextBlob:
 - <https://textblob.readthedocs.io/en/dev/quickstart.html#sentiment-analysis>
- **Thematic Analysis (Keyword Extraction & Topic Modeling):**
 - spaCy (for tokenization, lemmatization, POS tagging):
<https://spacy.io/usage/spacy-101>
 - Scikit-learn (for TF-IDF, LDA, NMF):
 - TF-IDF (TfidfVectorizer):
https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html
 - LDA (LatentDirichletAllocation):
<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.LatentDirichletAllocation.html>
 - NMF (NMF):
<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.NMF.html>
 - Hugging Face transformers (for Zero-Shot Classification - e.g., facebook/bart-large-mnli):
 - ZeroShotClassificationPipeline:
https://huggingface.co/docs/transformers/main_classes/pipelines#transformers.ZeroShotClassificationPipeline
 - Sentence Transformers (for embeddings if doing embedding-based clustering):
<https://www.sbert.net/docs/quickstart.html>
- **Database (PostgreSQL):**
 - PostgreSQL Official Documentation:
 - Getting Started: <https://www.postgresql.org/docs/current/tutorial-start.html>

- PostgreSQL Downloads: <https://www.postgresql.org/download/>
 - psycopg2 (Python adapter for PostgreSQL): Documentation: <https://www.psycopg.org/docs/>
 - Installation: <https://www.psycopg.org/docs/install.html>
 - PostgreSQL Tutorial: <https://docs.sqlalchemy.org/en/20/dialects/postgresql.html>
 - pgAdmin (GUI Tool for PostgreSQL): <https://www.pgadmin.org/>
 - Basic SQL Tutorials: [W3Schools SQL Tutorial:](https://www.w3schools.com/sql/) <https://www.w3schools.com/sql/>
 - PostgreSQL Tutorial: <https://www.postgresqltutorial.com/>
- **Unit Testing:**
 - <https://docs.pytest.org/en/stable/getting-started.html>
 - <https://docs.python.org/3/library/unittest.html#basic-example>