Revised Data Engineering Curriculum

Month 1: Foundation Building

Week 1: Python for Data Engineering

1. Quick Wins with Python:

   o Learn Python basics (data types, loops, and conditionals) by building small programs, such as a CSV data analyzer.

   o Hands-on Challenge: Write a script to parse a log file and extract useful information.

2. Data Manipulation & Visualization:

   o Use pandas to clean and manipulate datasets.

   o Visualize results with matplotlib or seaborn (e.g., show the top 5 items sold from a dataset).

   o Mini-Project: Create a script to process sales data and output a report with visualizations.

3. Automate Simple Tasks with Python:

   o Write Python scripts to automate common tasks like renaming files or scraping data from websites.

   o Hands-on Challenge: Write a program that automatically downloads and organizes files from a given URL.

Week 2: SQL

1. Foundations of SQL:

   o Learn the basics of SQL by querying a mock e-commerce database.

   o Interactive Challenge: Retrieve top-selling products, identify the highest-paying customers, and calculate monthly revenue.

2. SQL for Data Transformation:

   o Perform data cleansing and transformation with SQL.

   o Hands-On Activity: Normalize a poorly designed database into a clean, efficient structure.

3. Mini-Project: Build Your First SQL Portfolio:

   o Design a small database for a real-world scenario (e.g., a library system or movie database). Write queries to extract meaningful insights.

Week 3: Linux/Bash Scripting

1. Master Linux Basics:

   o Learn the Linux file system and basic commands (ls, grep, find, etc.).

   o Hands-On Activity: Automate file organization using bash scripting.

2. Write Bash Scripts for Data Tasks:

   o Write scripts to automate file cleaning, data merging, and scheduling tasks with cron jobs.

   o Challenge: Process and clean a folder of messy log files using bash.

3. Mini-Project:

   o Automate a simple ETL pipeline using Bash (e.g., extract files from a folder, clean them, and move them to an organized directory).

---

Week 4: Data Engineering Basics

1. Introduction to Data Pipelines:

   o Learn the concept of ETL (Extract, Transform, Load) and its importance in data engineering.

   o Hands-On Activity: Write a Python script to perform a basic ETL process on a small dataset.

2. Data Formats and Storage Systems:

   o Work with data formats like CSV, JSON, and Parquet.

   o Challenge: Convert a dataset into multiple formats and compress it.

3. Mini-Project:

   o Build an end-to-end pipeline to clean a dataset and load it into an SQLite database.

---

Month 2: Advanced Tools and Concepts

Week 5: Big Data Tools

1. Apache Spark Introduction:

   o Learn Spark fundamentals, including RDDs and DataFrames.

   o Hands-On Challenge: Analyze a large dataset with PySpark and calculate insights like average sales per region.

2. Spark SQL for Big Data:

   o Write SQL queries in Spark to process structured data.

- o Activity: Query and process data stored in a Parquet file.

3. Mini-Project:

   - o Process a large dataset using PySpark and generate a summary report.

---

Week 6: Workflow Orchestration

1. Apache Airflow Basics:

   - o Learn to create and deploy DAGs for task orchestration.

   - o Hands-On Activity: Schedule a daily task to clean and archive logs.

2. Advanced Airflow Features:

   - o Use XComs for data sharing between tasks, and integrate Airflow with cloud storage.

   - o Challenge: Orchestrate a multi-step pipeline involving data extraction, cleaning, and storage.

3. Mini-Project:

   - o Build and deploy an Airflow pipeline that ingests API data and loads it into a database.

---

Week 7: Data Modeling and Advanced SQL

1. Dimensional Data Modeling:

   - o Design star and snowflake schemas for business scenarios.

   - o Activity: Create a star schema for a sales reporting database.

2. Advanced SQL Techniques:

   - o Learn window functions, CTEs, and recursive queries.

   - o Challenge: Write complex queries to calculate rolling averages and rank data.

3. Mini-Project:

   - o Design a data warehouse schema and populate it with mock data.

---

Week 8: Azure Cloud for Data Engineering

1. Azure Data Storage and Processing:

   - o Learn Azure Blob Storage and Azure Data Factory basics.

   - o Activity: Store and retrieve files in Azure Blob Storage.

2. ETL with Azure Data Factory:

   o Build pipelines to process data in Azure.

   o Challenge: Set up an Azure Data Factory pipeline to clean data and save it in Blob Storage.

3. Mini-Project:

   o Build a cloud-based ETL pipeline using Azure Data Factory and Blob Storage.

---

Month 3: Real-World Applications

Week 9: Databricks for Big Data Analytics

1. Introduction to Databricks:

   o Overview of Databricks and its role in big data processing.

   o Setting up a Databricks workspace and understanding clusters, notebooks, and jobs.

2. Mounting Azure Data Lake on Databricks:

   o Connect Databricks to Azure Data Lake and access large datasets.

   o Hands-on Activity: Mount an Azure Data Lake container to your Databricks workspace.

3. Analyzing Large Datasets:

   o Process and analyze large datasets using Spark DataFrames in Databricks.

   o Perform transformations, aggregations, and visualizations in Databricks notebooks.

   o Challenge: Analyze a 1GB+ dataset stored in Azure Data Lake and generate key business insights.

4. Mini-Project:

   o Build an end-to-end data analysis solution:

      ▪ Mount an Azure Data Lake container on Databricks.

      ▪ Extract, transform, and analyze data using Spark.

      ▪ Visualize the results in Databricks notebooks.

---

Week 10: DevOps for Data Engineering

1. CI/CD Pipelines:

   o Set up a CI/CD pipeline for deploying data pipelines.

      o   Activity: Use GitHub Actions to test and deploy a data pipeline.

2. Docker for Data Engineering:

      o   Learn to containerize data engineering workflows.

      o   Challenge: Deploy a pipeline in Docker and scale it.

3. Mini-Project:

      o   Build and deploy a containerized ETL pipeline with a CI/CD process.

---

Week 11–12: Capstone Project

1. Design an End-to-End Solution:

      o   Create a complete ETL pipeline:

            ▪   Extract data from APIs, files, or databases.

            ▪   Transform data with Python/Spark.

            ▪   Load data into an Azure SQL database or a data warehouse.

2. Enhance with Streaming and Cloud:

      o   Add a real-time streaming component with Kafka and Spark.

      o   Deploy the solution to Azure with containerization and CI/CD.

3. Portfolio Presentation:

      o   Document and present the solution as a portfolio project.