

# Data Engineering Nanodegree Program Syllabus



*Build Production Ready Data Warehouses at Scale*

---

## Before You Start

**Prerequisites:** Students should have intermediate SQL and Python programming skills.

**Educational Objectives:** Students will learn to

- Create user-friendly relational and NoSQL data models
- Create scalable and efficient data warehouses
- Identify the appropriate use cases for different big data technologies
- Work efficiently with massive datasets
- Build and interact with a cloud-based data lake
- Automate and monitor data pipelines
- Develop proficiency in Spark, Airflow, and AWS tools

**Estimated Length of Program:** 5 months

**Program Structure:** Self-paced

**Textbooks Required:** None

# Course 1: Data Modeling

In this course, you'll learn to create relational and NoSQL data models to fit the diverse needs of data consumers. You'll understand the differences between different data models, and how to choose the appropriate data model for a given situation. You'll also build fluency in PostgreSQL and Apache Cassandra.

Lesson Title	Learning Outcomes
<b>Introduction to Data Modeling</b>	<ul style="list-style-type: none"><li>→ Understand the purpose of data modeling</li><li>→ Identify the strengths and weaknesses of different types of databases and data storage techniques</li><li>→ Create a table in Postgres and Apache Cassandra</li></ul>
<b>Relational Data Models</b>	<ul style="list-style-type: none"><li>→ Understand when to use a relational database</li><li>→ Understand the difference between OLAP and OLTP databases</li><li>→ Create normalized data tables</li><li>→ Implement denormalized schemas (e.g. STAR, Snowflake)</li></ul>
<b>NoSQL Data Models</b>	<ul style="list-style-type: none"><li>→ Understand when to use NoSQL databases and how they differ from relational databases</li><li>→ Select the appropriate primary key and clustering columns for a given use case</li><li>→ Create a NoSQL database in Apache Cassandra</li></ul>

## Projects: Data Modeling with Postgres and Apache Cassandra

In these projects, you'll model user activity data for a music streaming app called Sparkify. You'll create a database and ETL pipeline, in both Postgres and Apache Cassandra, designed to optimize queries for understanding what songs users are listening to. For PostgreSQL, you will also define Fact and Dimension tables and insert data into your new tables. For Apache Cassandra, you will model your data so you can run specific queries provided by the analytics team at Sparkify.

## Course 2: Cloud Data Warehouses

In this course, you'll learn to create cloud-based data warehouses. You'll sharpen your data warehousing skills, deepen your understanding of data infrastructure, and be introduced to data engineering on the cloud using Amazon Web Services (AWS).

Lesson Title	Learning Outcomes
<b>Introduction to the Data Warehouses</b>	<ul style="list-style-type: none"><li>→ Understand Data Warehousing architecture</li><li>→ Run an ETL process to denormalize a database (3NF to Star)</li><li>→ Create an OLAP cube from facts and dimensions</li><li>→ Compare columnar vs. row oriented approaches</li></ul>
<b>Introduction to the Cloud with AWS</b>	<ul style="list-style-type: none"><li>→ Understand cloud computing</li><li>→ Create an AWS account and understand their services</li><li>→ Set up Amazon S3, IAM, VPC, EC2, RDS PostgreSQL</li></ul>
<b>Implementing Data Warehouses on AWS</b>	<ul style="list-style-type: none"><li>→ Identify components of the Redshift architecture</li><li>→ Run ETL process to extract data from S3 into Redshift</li><li>→ Set up AWS infrastructure using Infrastructure as Code (IaC)</li><li>→ Design an optimized table by selecting the appropriate distribution style and sorting key</li></ul>

### Project 2: Data Infrastructure on the Cloud

In this project, you are tasked with building an ELT pipeline that extracts their data from S3, stages them in Redshift, and transforms data into a set of dimensional tables for their analytics team to continue finding insights in what songs their users are listening to.

## Course 3: Data Lakes with Spark

In this course, you will learn more about the big data ecosystem and how to use Spark to work with massive datasets. You'll also learn about how to store big data in a data lake and query it with Spark.

Lesson Title	Learning Outcomes
<b>The Power of Spark</b>	<ul style="list-style-type: none"><li>→ Understand the big data ecosystem</li><li>→ Understand when to use Spark and when not to use it</li></ul>
<b>Data Wrangling with Spark</b>	<ul style="list-style-type: none"><li>→ Manipulate data with SparkSQL and Spark Dataframes</li><li>→ Use Spark for ETL purposes</li></ul>
<b>Debugging and Optimization</b>	<ul style="list-style-type: none"><li>→ Troubleshoot common errors and optimize their code using the Spark WebUI</li></ul>
<b>Introduction to Data Lakes</b>	<ul style="list-style-type: none"><li>→ Understand the purpose and evolution of data lakes</li><li>→ Implement data lakes on Amazon S3, EMR, Athena, and Amazon Glue</li><li>→ Use Spark to run ELT processes and analytics on data of diverse sources, structures, and vintages</li><li>→ Understand the components and issues of data lakes</li></ul>

### Project 3: Big Data with Spark

In this project, you'll build an ETL pipeline for a data lake. The data resides in S3, in a directory of JSON logs on user activity on the app, as well as a directory with JSON metadata on the songs in the app. You will load data from S3, process the data into analytics tables using Spark, and load them back into S3. You'll deploy this Spark process on a cluster using AWS.

## Project 4: Automate Data Pipelines

In this course, you'll learn to schedule, automate, and monitor data pipelines using Apache Airflow. You'll learn to run data quality checks, track data lineage, and work with data pipelines in production.

Lesson Title	Learning Outcomes
<b>Data Pipelines</b>	<ul style="list-style-type: none"><li>→ Create data pipelines with Apache Airflow</li><li>→ Set up task dependencies</li><li>→ Create data connections using hooks</li></ul>
<b>Data Quality</b>	<ul style="list-style-type: none"><li>→ Track data lineage</li><li>→ Set up data pipeline schedules</li><li>→ Partition data to optimize pipelines</li><li>→ Write tests to ensure data quality</li><li>→ Backfill data</li></ul>
<b>Production Data Pipelines</b>	<ul style="list-style-type: none"><li>→ Build reusable and maintainable pipelines</li><li>→ Build your own Apache Airflow plugins</li><li>→ Implement subDAGs</li><li>→ Set up task boundaries</li><li>→ Monitor data pipelines</li></ul>

### Project: Data Pipelines with Airflow

In this project, you'll continue your work on the music streaming company's data infrastructure by creating and automating a set of data pipelines. You'll configure and schedule data pipelines with Airflow and monitor and debug production pipelines.

## Data Engineering Nanodegree Capstone Project

The purpose of the data engineering capstone project is to give you a chance to combine what you've learned throughout the program. This project will be an important part of your portfolio that will help you achieve your data engineering-related career goals.

In this project, you'll define the scope of the project and the data you'll be working with. We'll provide guidelines, suggestions, tips, and resources to help you be successful, but your project will be unique to you. You'll gather data from several different data sources; transform, combine, and summarize it; and create a clean database for others to analyze.