



# Data Engineer Nanodegree Syllabus

Build Production Ready Data Warehouses at Scale



UDACITY  
FOR ENTERPRISE



## Table of Contents

Before You Start	3
Contact Info	3
Nanodegree Program Info	3
Project 1: Data Modeling with Postgress and Apache Cassandra	4
Project 2: Data Infrastructure on the Cloud	5
Project 3: Big Data with Spark	5
Project 4: Data Pipelines with Airflow	6
Project 5: Data Engineering Nanodegree Capstone Project	7



## 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

## Contact Info

While going through the program, if you have questions about anything, you can reach us at [enterprise-support@udacity.com](mailto:enterprise-support@udacity.com). For help from Udacity mentors and peers, please visit the Udacity classroom.

## Nanodegree Program Info

### Technical Requirements

**Hardware Requirements:** webcam, microphone

**Software and Software Version Requirements:** No software required

**LENGTH OF PROGRAM\*:** 6 months

**FREQUENCY OF CLASSES:** Self-paced

**TEXTBOOKS REQUIRED:** None

\*This is a self-paced program and the length is an estimation of total hours the average student may take to complete all required coursework, including lecture and project time. Actual hours may vary.



## Project 1: 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.

### Supporting Lesson Content: Data Modeling

Lesson Title	Learning Outcomes
INTRODUCTION TO DATA MODELING	<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>
RELATIONAL DATA MODELS	<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>
NOSQL DATA MODELS	<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>

## Project 2: Data Infrastructure on the Cloud

In this project, you are tasked with building an ETL 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.

### Supporting Lesson Content: Cloud Data Warehouses

Lesson Title	Learning Outcomes
INTRODUCTION TO THE DATA WAREHOUSES	<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>
INTRODUCTION TO THE CLOUD WITH AWS	<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>
IMPLEMENTING DATA WAREHOUSES ON AWS	<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 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.

## Supporting Lesson Content: Data Lakes with Spark

Lesson Title	Learning Outcomes
THE POWER OF SPARK	<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>
DATA WRANGLING WITH SPARK	<ul style="list-style-type: none"><li>• Manipulate data with SparkSQL and Spark Dataframes.</li><li>• Use Spark for ETL purposes.</li></ul>
DEBUGGING AND OPTIMIZATION	<ul style="list-style-type: none"><li>• Troubleshoot common errors and optimize their code using the Spark WebUI.</li></ul>
INTRODUCTION TO DATA LAKES	<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 4: 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.

## Supporting Lesson Content: Automate Data Pipelines

Lesson Title	Learning Outcomes
DATA PIPELINES	<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>

## Supporting Lesson Content: Data Lakes with Spark (Continued)

Lesson Title	Learning Outcomes
DATA QUALITY	<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>
PRODUCTION DATA PIPELINES	<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 5: 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.



Learn more at [www.udacity.com/enterprise](http://www.udacity.com/enterprise)