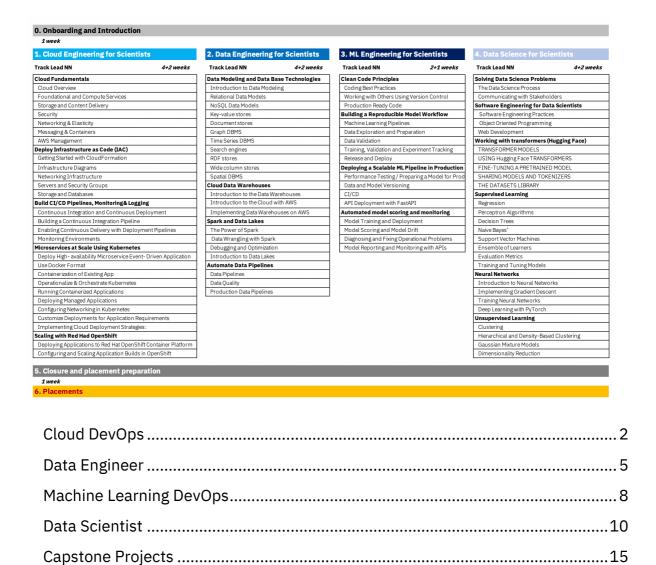
Software Engineering for Scientist

Extended curriculum overview (subject to changes)

This document provides an overview on the content and learning objectives from the different streams as well as the capstone projects.

The contents are subject to changes introduced by the steam leads.

As far as possible and for the pilot phase, the program will rely on existing content from renowned platforms (Udacity, Coursera, etc).



Cloud DevOps

Module	Learning Objectives
Cloud Fundamentals	
Cloud Overview	Learn the basics of cloud computing including cloud deployment models, benefits, and popular options Explore services provided by Amazon Web Services(AWS)
Foundational and Compute Services	Learn why we need servers, compute power, and security Explore AWS compute services like Elastic Cloud Compute (EC2), Virtual Private Cloud (VPC), Lambda for serverless framework, and Elastic Beanstalk in action Launch a secure EC2 instance, create and execute a Lambda, and deploy an application to Elastic Beanstalk
Storage and Content Delivery	 Learn why we need storage and content delivery in the cloud Learn storage services like S3, DynamoDB, Relational Database Service (RDS), and CloudFront Create a DynamoDB table, launch a MySQL database instance, and create a CloudFront distribution
Security	Learn the importance of security in the cloud See Identity & Access Management (IAM) in action Secure applications using IAM users, groups, and policies
Networking & Elasticity	 Learn the basics of networking and elasticity in the cloud Examine services like Route 53, EC2 Auto Scaling, and Elastic Load Balancing Add an auto scaling policy to your EC2 instance
Messaging & Containers	 Learn the basics of messaging and containers in the cloud Explore services like Simple Notification Service (SNS), Simple Queue Service (SQS), and Elastic Container Service (ECS) Create cloud notifications using SNS
AWS Management	Learn why we need logging, auditing, and resource management in the cloud Understand services like CloudWatch, CloudTrail, CloudFormation, and the AWS Command Line Interface (CLI) Explore the CLI
Deploy Infrastructure	as Code (IAC)
Getting Started with CloudFormation Infrastructure	Set up the necessary tools to get started with CloudFormation and deploy your first server using CloudFormation Convert business requirements into infrastructure diagrams and understand the principles behind design chaines.
Diagrams Networking Infrastructure	Principles behind design choices Implement a virtual private network and subnets and learn how to provide inbound and outbound internet access to your public and private subnets inside your VPC Use routing table to route the traffic within your virtual private cloud

Servers and	Deploy a web server into an autoscaling group
Security Groups	Implement load-balancer to increase capacity of your app
cocarri, arreape	Implement security groups and understand the concept of
	least-privilege as it applies to network traffic
	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
Storage and	Deploy S3 storage for images, config files, and more.
Databases	Deploy relational database and encryption service for your application
Build CI/CD Pipelines	 , Monitoring & Logging
Continuous	Understand the fundamentals of CI/CD.
	, ,
Integration and	Give examples of business-centered benefits of CI/CD. Type in a the utility of centing and alivery in a devite and a second control of the control
Continuous	• Examine the utility of continuous delivery in a dev team.
Deployment	• List best practices.
	Differentiate deployment strategies.
	Recognize common building blocks of CI/CD pipelines.
Building a	Understand how and why to use configuration management tools.
Continuous	Utilize a configuration management tool to accomplish deployment to cloud-based
Integration Pipeline	servers.
integration ripetine	Design a complete CI pipeline.
 Enabling	Know what configuration management tools are and how to use them.
Continuous Delivery	Design an Ansible Playbook and control a remote machine.
with Deployment	Build an Ansible Inventory file.
Pipelines	Make various types of CD jobs in a CI/CD pipeline.
ripetities	• Make various types of CD Jobs III a CI/CD pipetifie.
Monitoring	• Install and configure Prometheus as a monitoring tool.
Environments	Get various data sources into Prometheus.
	Analyze monitoring data.
	Set up alerts.
Microservices at Scale	Using Kubernetes
Deploy High-	Understand Serverless (AWS Lambda) concepts
availability	Understand which container abstraction to use: AWS Lambda or Kubernetes
Microservice Event-	Deploy producer/consumer AWS Lambda applications
Driven Application	Configure CloudWatch events
	Company crosses
Use Docker Format	Understand Docker image format Containers
	Run and modify Docker containers locally
	Deploy customized containers to Amazon ECR
Containerization of	Use the appropriate Docker base image
Existing App	• Install packages into Docker image
	Copy application into Docker image
	Configure application setup and start in Docker image

Operationalize &	Understand Kubernetes concepts
Orchestrate	Configure monitoring, alerts, and incidence response
Kubernetes	Integrate CI/CD Pipeline
	Configure Autoscaling
	3 3 3 3 3 3 3 3 3 3
Running	Contrasting Kubernetes Distributions
Containerized	Introducing Kubectl
Applications	Running and Interacting with Your First Application
Deploying Managed	Managing Containers
Applications	Creating a Deployment
	Understanding the Schema of a Deployment Resource
Configuring	Kubernetes Networking
Networking in	• Introducing Kubernetes Services
Kubernetes	Discovering Kubernetes Services
	Kubernetes Ingress + Kubernetes Ingress Controller
	• Ingress Resource Configuration + Testing Your Ingress
Customize	Defining Resource Requests and Limits for Pods + Viewing Requests, Limits, and
Deployments for	Actual Usage
Application	Applying Quotas and Limit Ranges
Requirements	Proving Liveness, Readiness and Startup
	Methods of Checking Application Health + Creating Probes
	Externalizing Application Configuration in Kubernetes + Using Secret and
	Configuration Map Resources
	Creating and Managing Secrets and Configuration Maps
	• Injecting Data from Secrets and Configuration Maps into Applications + Application
	Configuration Options
	Exploring Environment Variables
Implementing	Deployment Strategies in Kubernetes
Cloud Deployment	Implementing Advanced Deployment Strategies Using the Kubernetes Router
Strategies:	Guided Exercise: Implementing Cloud Deployment Strategies
Scaling with Red Had	OpenShift
5 1 :	
Deploying	Introducing OpenShift Container Platform The desired Con
Applications to Red	Introducing the Developer Web Console for OpenShift
Hat OpenShift	
Container Platform	
Configuring and	Updating an Application Configuration Application Country
Scaling Application	Configuring Application Secrets Connecting on Application to a Database
Builds in OpenShift	Connecting an Application to a Database Conline Applications in One of Shift
	Scaling Applications in OpenShift

Data Engineer

Identify the strengths and weaknesses of different types of databases and data storage techniques	Module	Learning Objectives
Identify the strengths and weaknesses of different types of databases and data storage techniques	Data Modeling and Da	ata Base Technologies
Identify the strengths and weaknesses of different types of databases and data storage techniques Create a table in Postgres and Apache Cassandra Understand when to use a relational database Understand the difference between OLAP and OLTP databases Create normalized data tables Implement denormalized schemas (e.g. STAR, Snowflake) VosQL Data Models Understand when to use NoSQL databases and how they differ from relational databases Select the appropriate primary key and clustering columns for a given use case Create a NoSQL database in Apache Cassandra Vunderstand when to use a key value store Overview of the main key-value store technologies Create a key-value store database in Redis Document stores Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Understand when to use a time series database Overview of the main time series database Coverview of the main time series database Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main min RDF store technologies	Introduction to Data	Understand the purpose of data modeling
of databases and data storage techniques • Create a table in Postgres and Apache Cassandra **Understand when to use a relational database • Understand the difference between OLAP and OLTP databases • Create normalized data tables • Implement denormalized schemas (e.g. STAR, Snowflake) **NoSQL Data Models** • Understand when to use NoSQL databases and how they differ from relational databases • Select the appropriate primary key and clustering columns for a given use case • Create a NoSQL database in Apache Cassandra **Key-value stores** • Understand when to use a key value store • Overview of the main key-value store technologies • Create a key-value store database in Redis **Document stores** • Understand when to use a document store • Overview of the main document store technologies • Create a document store database in MongoDB **Graph DBMS** • Understand when to use a graph database • Overview of the main graph database technologies • Create a graph database in Neo4j (or ArangoDB) **Time Series DBMS** • Understand when to use a time series database • Overview of the main time series database technologies • Create a time series database in InfluxDB (o TimescaleDB) **Search engines** • Understand when to use a search engine • Overview of the main search engine technologies • Create a search engine using Elasticsearch **RDF stores** • Understand when to use a RDF store • Overview of the main RDF store technologies • Create a RDF store using Virtuoso (o MarkLogic) **Wide column stores** • Understand when to use a wide column database • Overview of the main wide column database • Overview of the main wide column database • Overview of the main mill response technologies	Modeling	
Create a table in Postgres and Apache Cassandra Understand when to use a relational database Understand the difference between OLAP and OLTP databases Create normalized data tables Implement denormalized schemas (e.g. STAR, Snowflake) NoSQL Data Models Understand when to use NoSQL databases and how they differ from relational databases Select the appropriate primary key and clustering columns for a given use case Create a NoSQL database in Apache Cassandra Key-value stores Understand when to use a key value store Overview of the main key-value store technologies Create a key-value store database in Redis Document stores Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Understand when to use a time series database Overview of the main time series database Overview of the main time series database Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine Overview of the main search engine technologies Create a search engine using Elasticsearch Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database Overview of the main wide column database technologies		
Understand the difference between OLAP and OLTP databases		
databases	Relational Data	Understand when to use a relational database
Create normalized data tables Implement denormalized schemas (e.g. STAR, Snowflake) ONOSQL Data Models Understand when to use NoSQL databases and how they differ from relational databases Select the appropriate primary key and clustering columns for a given use case Create a NoSQL database in Apache Cassandra Understand when to use a key value store Overview of the main key-value store technologies Create a key-value store database in Redis Obcument stores Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Time Series DBMS Understand when to use a time series database Overview of the main time series database Overview of the main time series database Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Understand when to use a wide column database Overview of the main wide column database technologies	Models	Understand the difference between OLAP and OLTP
Implement denormalized schemas (e.g. STAR, Snowflake) Understand when to use NoSQL databases and how they differ from relational databases Select the appropriate primary key and clustering columns for a given use case Create a NoSQL database in Apache Cassandra Understand when to use a key value store Overview of the main key-value store technologies Create a key-value store database in Redis Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Time Series DBMS Understand when to use a time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Understand when to use a wide column database Overview of the main wide column database Overview of the main wide column database		databases
NoSQL Data Models - Understand when to use NoSQL databases and how they differ from relational databases - Select the appropriate primary key and clustering columns for a given use case - Create a NoSQL database in Apache Cassandra - Understand when to use a key value store - Overview of the main key-value store technologies - Create a key-value store database in Redis - Understand when to use a document store - Overview of the main document store technologies - Create a document store database in MongoDB - Understand when to use a graph database - Overview of the main graph database - Overview of the main graph database - Create a graph database in Neo4j (or ArangoDB) - Understand when to use a time series database - Overview of the main time series database technologies - Create a time series database in InfluxDB (o TimescaleDB) - Understand when to use a search engine - Overview of the main search engine technologies - Create a search engine using Elasticsearch - Understand when to use a RDF store - Overview of the main RDF store technologies - Create a RDF store using Virtuoso (o MarkLogic) - Understand when to use a wide column database - Overview of the main wide column database		Create normalized data tables
databases		• Implement denormalized schemas (e.g. STAR, Snowflake)
databases	NoSOL Data Models	Understand when to use NoSOL databases and how they differ from relational
Create a NoSQL database in Apache Cassandra Founderstand when to use a key value store Overview of the main key-value store technologies Create a key-value store database in Redis Document stores Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Time Series DBMS Understand when to use a time series database Overview of the main time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database Overview of the main wide column database	•	•
Create a NoSQL database in Apache Cassandra Founderstand when to use a key value store Overview of the main key-value store technologies Create a key-value store database in Redis Document stores Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Time Series DBMS Understand when to use a time series database Overview of the main time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database Overview of the main wide column database		Select the appropriate primary key and clustering columns for a given use case
Overview of the main key-value store technologies Create a key-value store database in Redis Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Time Series DBMS Understand when to use a time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Understand when to use a wide column database Overview of the main wide column database Overview of the main wide column database		
Create a key-value store database in Redis Understand when to use a document store Overview of the main document store technologies Create a document store database in MongoDB Graph DBMS Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Time Series DBMS Understand when to use a time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Understand when to use a wide column database Overview of the main wide column database	Key-value stores	Understand when to use a key value store
Document stores - Understand when to use a document store - Overview of the main document store technologies - Create a document store database in MongoDB - Understand when to use a graph database - Overview of the main graph database technologies - Create a graph database in Neo4j (or ArangoDB) - Understand when to use a time series database - Overview of the main time series database - Overview of the main time series database technologies - Create a time series database in InfluxDB (o TimescaleDB) - Understand when to use a search engine - Overview of the main search engine technologies - Create a search engine using Elasticsearch - Understand when to use a RDF store - Overview of the main RDF store technologies - Create a RDF store using Virtuoso (o MarkLogic) - Understand when to use a wide column database - Overview of the main wide column database technologies		Overview of the main key-value store technologies
Overview of the main document store technologies Create a document store database in MongoDB Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Understand when to use a time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Understand when to use a wide column database Overview of the main wide column database		Create a key-value store database in Redis
Create a document store database in MongoDB Understand when to use a graph database Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Understand when to use a time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies	Document stores	Understand when to use a document store
Graph DBMS - Understand when to use a graph database - Overview of the main graph database technologies - Create a graph database in Neo4j (or ArangoDB) Time Series DBMS - Understand when to use a time series database - Overview of the main time series database technologies - Create a time series database in InfluxDB (o TimescaleDB) Search engines - Understand when to use a search engine - Overview of the main search engine technologies - Create a search engine using Elasticsearch RDF stores - Understand when to use a RDF store - Overview of the main RDF store technologies - Create a RDF store using Virtuoso (o MarkLogic) Wide column stores - Understand when to use a wide column database - Overview of the main wide column database technologies		Overview of the main document store technologies
Overview of the main graph database technologies Create a graph database in Neo4j (or ArangoDB) Understand when to use a time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies		Create a document store database in MongoDB
Create a graph database in Neo4j (or ArangoDB) Understand when to use a time series database Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Search engines Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies	Graph DBMS	• Understand when to use a graph database
Time Series DBMS - Understand when to use a time series database - Overview of the main time series database technologies - Create a time series database in InfluxDB (o TimescaleDB) Search engines - Understand when to use a search engine - Overview of the main search engine technologies - Create a search engine using Elasticsearch RDF stores - Understand when to use a RDF store - Overview of the main RDF store technologies - Create a RDF store using Virtuoso (o MarkLogic) Wide column stores - Understand when to use a wide column database - Overview of the main wide column database technologies		Overview of the main graph database technologies
Overview of the main time series database technologies Create a time series database in InfluxDB (o TimescaleDB) Ounderstand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch Ouderstand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Ouderstand when to use a wide column database Overview of the main wide column database technologies		Create a graph database in Neo4j (or ArangoDB)
Create a time series database in InfluxDB (o TimescaleDB) Understand when to use a search engine Overview of the main search engine technologies Create a search engine using Elasticsearch RDF stores Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies	Time Series DBMS	Understand when to use a time series database
Search engines - Understand when to use a search engine - Overview of the main search engine technologies - Create a search engine using Elasticsearch RDF stores - Understand when to use a RDF store - Overview of the main RDF store technologies - Create a RDF store using Virtuoso (o MarkLogic) Wide column stores - Understand when to use a wide column database - Overview of the main wide column database technologies		
Overview of the main search engine technologies Create a search engine using Elasticsearch Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies		Create a time series database in InfluxDB (o TimescaleDB)
Create a search engine using Elasticsearch Understand when to use a RDF store Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies	Search engines	Understand when to use a search engine
RDF stores • Understand when to use a RDF store • Overview of the main RDF store technologies • Create a RDF store using Virtuoso (o MarkLogic) Wide column stores • Understand when to use a wide column database • Overview of the main wide column database technologies		Overview of the main search engine technologies
Overview of the main RDF store technologies Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies		Create a search engine using Elasticsearch
Create a RDF store using Virtuoso (o MarkLogic) Wide column stores Understand when to use a wide column database Overview of the main wide column database technologies	RDF stores	
Wide column stores • Understand when to use a wide column database • Overview of the main wide column database technologies		_
Overview of the main wide column database technologies		Create a RDF store using Virtuoso (o MarkLogic)
	Wide column stores	
Create a wide column database using Cassandra (or HBase)		
		Create a wide column database using Cassandra (or HBase)

Spatial DBMS	Understand when to use a Spatial DBMS
·	Overview of the main Spatial DBMS technologies
	Create a Spatial DBMS using PostGIS
Cloud Data Warehous	ses
Introduction to the	Understand Data Warehousing architecture
Data Warehouses	• Run an ETL process to denormalize a database (3NF to Star)
	Create an OLAP cube from facts and dimensions
	Compare columnar vs. row oriented approaches
Introduction to the	Understand cloud computing
Cloud with AWS	Create an AWS account and understand their services
(refresher)	Set up Amazon S3, IAM, VPC, EC2, RDS PostgreSQL
Implementing Data	Identify components of the Redshift architecture
Warehouses on AWS	• Run ETL process to extract data from S3 into Redshift
	Set up AWS infrastructure using Infrastructure as Code (IaC) Design an antimized table by selecting the appropriate distribution style and
	• Design an optimized table by selecting the appropriate distribution style and sorting key
Spark and Data Lakes	
The Power of Spark	Understand the big data ecosystem
	Understand when to use Spark and when not to use it
Data Wrangling with	Manipulate data with SparkSQL and Spark Dataframes
Spark	Use Spark for ETL purposes
Debugging and Optimization	Troubleshoot common errors and optimize their code using the Spark WebUI
Introduction to Data	Understand the purpose and evolution of data lakes
Lakes	• Implement data lakes on Amazon S3, EMR, Athena, and Amazon Glue
	Use Spark to run ELT processes and analytics on data of diverse sources,
	structures, and vintagesUnderstand the components and issues of data lakes
Automate Data Pipeli	· · · · · · · · · · · · · · · · · · ·
Data Pipelines	Create data pipelines with Apache Airflow
•	• Set up task dependencies
	Create data connections using hooks
Data Quality	Track data lineage
	Set up data pipeline schedules
	Partition data to optimize pipelines
	Write tests to ensure data quality Backfill data
Production Data	Build reusable and maintainable pipelines
Pipelines	Build your own Apache Airflow plugins
	Implement subDAGs
	• Set up task boundaries
	Monitor data pipelines

Machine Learning DevOps

Module	Learning Objectives
Clean Code Principles	
Coding Best Practices	Write clean, modular and well-documented code
_	Refactor code for efficiency
	• Follow PEP8 Standards
	Automate use of PEP8 standards using PyLint and AutoPEP8
Working with Others	Work independently using git and Github
Using Version Control	Work with teams using git and Github
	Create branches for isolating changes in git and Github
	Open pull requests for making changes to production code
	Conduct and receive code reviews using best practices
Production Ready	Correctly use try-except blocks to identify errors
Code	Create unit tests to test programs
	Track actions and results of processes with logging
	Identify model drift and when automated or non-automated retraining should be
	used to make model updates
Building a Reproducib	le Model Workflow
Machine Learning	MLOps fundamentals
Pipelines	Version data and artifacts
	Write a ML pipeline component
	Link together ML components
Data Exploration and	Execute and track the Exploratory Data Analysis(EDA)
Preparation	Clean and pre-process the data
	Segregate(split)datasets
Data Validation	Use pytest with parameters for reproducible and automatic data tests
	Perform deterministic and non-deterministic data tests
Training, Validation	Tame the chaos with experiment, code and data tracking
and Experiment	Track experiments with W&B
Tracking	Validate and choose best-performing model
	Export model as an inference artifact
	Test final inference artifact
Release and Deploy	Release pipeline code
	Options for deployment and how to deploy a model
Deploying a Scalable N	ML Pipeline in Production

Performance Testing and Preparing a Model for Production	 Analyze slices of data when training and testing models Probe a model for bias using common frameworks such as Aequitas Write model cards that explain the purpose, provenance and pitfalls of a model
Data and Model	Version control data/models/etc locally using DVC
Versioning	Setup remote storage for use with DVC
	Create pipelines and track experiments with DVC
CI/CD	Follow software engineering principles by automating, testing and versioning code
	Setup Continuous Integration using GitHub Actions
	Setup Continuous Deployment using Heroku
API Deployment with	Write an API for machine learning inference using FastAPI
FastAPI	Deploy a machine learning inference API to Heroku
	Write unit tests for APIs using the requests module
Automated model sco	ring and monitoring
Model Training and	• Ingest data
Deployment	Automatically train models
	Deploy models to production
	Keep records about processes
	Automate processes using cronjobs
Model Scoring and	Automatically score ML models
Model Drift	Keep records of model scores
	Check for model drift using several different model drift tests
	Determine whether models need to be retrained and re-deployed
Diagnosing and Fixing	Check data integrity and stability
Operational Problems	Check for dependency issues
	Check for timing issues
	Resolve operational issues
Model Reporting and	Create API endpoints that enable users to access model results, metrics and
Monitoring with APIs	diagnostics
Trontoring With 711 13	Set up APIs with multiple, complex endpoints
	Call APIs and work with their results
	Odit/ii 10 dila Work With their redute

Data Scientist

Module	Learning Objectives
Solving Data Science P	roblems
The Data Science	Apply the CRISP-DM process to business applications
Process	Wrangle, explore, and analyze a dataset
1100033	Apply machine learning for prediction
	Apply statistics for descriptive and inferential understanding
	Draw conclusions that motivate others to act on your results
	,
Communicating with	Implement best practices in sharing your code and written summaries
Stakeholders	Learn what makes a great data science blog
	• Learn how to create your ideas with the data science
	community
Software Engineering f	or Data Scientists
Software Engineering	Write clean, modular, and well-documented code
Practices	Refactor code for efficiency
	Create unit tests to test programs
	Write useful programs in multiple scripts
	Track actions and results of processes with logging
	Conduct and receive code reviews
Object Oriented	Understand when to use object oriented programming
Programming	Build and use classes
1 108.0	Understand magic methods
	Write programs that include multiple classes, and follow
	good code structure
	Learn how large, modular Python packages, such as pandas
	and scikit-learn, use object oriented programming • Portfolio Exercise: Build your
	own Python package
Web Development	Learn about the components of a web app
web bevelopment	Build a web application that uses Flask, Plotly, and the
	Bootstrap framework
	Portfolio Exercise: Build a data dashboard using a dataset
	of your choice and deploy it to a web application
	or your choice and deploy it to a web application
Data Engineering for Da	ı ata Scientists
_	

ETL Pipelines	 Understand what ETL pipelines are Access and combine data from CSV, JSON, logs, APIs, and databases Standardize encodings and columns Normalize data and create dummy variables Handle outliers, missing values, and duplicated data Engineer new features by running calculations Build a SQLite database to store cleaned data
Natural Language Processing	 Prepare text data for analysis with tokenization, lemmatization, and removing stop words Use scikit-learn to transform and vectorize text data Build features with bag of words and tf-idf Extract features with tools such as named entity recognition and part of speech tagging Build an NLP model to perform sentiment analysis
Machine Learning Pipelines	Understand the advantages of using machine learning pipelines to streamline the data preparation and modeling process Chain data transformations and an estimator with scikit- learn's Pipeline Use feature unions to perform steps in parallel and create more complex workflows Grid search over pipeline to optimize parameters for entire workflow Complete a case study to build a full machine learning pipeline that prepares data and creates a model for a dataset
Experiment Design and	Recommendations
Experiment Design	Understand how to set up an experiment, and the ideas associated with experiments vs. observational studies Defining control and test conditions • Choosing control and testing groups
Statistical Concerns of Experimentation	 Applications of statistics in the real world Establishing key metrics SMART experiments: Specific, Measurable, Actionable, Realistic, Timely
A/B Testing	 How it works and its limitations Sources of Bias: Novelty and Recency Effects Multiple Comparison Techniques (FDR, Bonferroni, Tukey) Portfolio Exercise: Using a technical screener from Starbucks to analyze the results of an experiment and write up your findings
Introduction to Recommendation Engines	 Distinguish between common techniques for creating recommendation engines including knowledge based, content based, and collaborative filtering based methods. Implement each of these techniques in python. List business goals associated with recommendation engines, and be able to recognize which of these goals are most easily met with existing

Matrix Factorization for Recommendations	 Understand the pitfalls of traditional methods and pitfalls of measuring the influence of recommendation engines under traditional regression and classification techniques. Create recommendation engines using matrix factorization and FunkSVD Interpret the results of matrix factorization to better understand latent features of customer data Determine common pitfalls of recommendation engines like the cold start problem and difficulties associated with usual tactics for assessing the effectiveness of recommendation engines using usual techniques, and potential solutions.
Data Science Projects	
Elective 1: Dog Breed Classification	 Use convolutional neural networks to classify different dogs according to their breeds Deploy your model to allow others to upload images of their dogs and send them back the corresponding breeds. Complete one of the most popular projects in Udacity history, and show the world how you can use your deep learning skills to entertain an audience!
Elective 2: Starbucks	 Use purchasing habits to arrive at discount measures to obtain and retain customers Identify groups of individuals that are most likely to be responsive to rebates.
Elective 3: Arvato Financial Services	Work through a real-world dataset and challenge provided by Arvato Financial Services, a Bertelsmann company Top performers have a chance at an interview with Arvato or another Bertelsmann company!
Elective 4: Spark for Big Data	Take a course on Apache Spark and complete a project using a massive, distributed dataset to predict customer churn Learn to deploy your Spark cluster on either AWS or IBM Cloud
Elective 5: Your Choice	Use your skills to tackle any other project of your choice
Data Science Projects	
Elective 1: Dog Breed Classification	 Use convolutional neural networks to classify different dogs according to their breeds Deploy your model to allow others to upload images of their dogs and send them back the corresponding breeds. Complete one of the most popular projects in Udacity history, and show the world how you can use your deep learning skills to entertain an audience!
Elective 2: Starbucks	 Use purchasing habits to arrive at discount measures to obtain and retain customers Identify groups of individuals that are most likely to be responsive to rebates.
Elective 3: Arvato Financial Services	 Work through a real-world dataset and challenge provided by Arvato Financial Services, a Bertelsmann company Top performers have a chance at an interview with Arvato or another Bertelsmann company!
Elective 4: Spark for Big Data	Take a course on Apache Spark and complete a project using a massive, distributed dataset to predict customer churn Learn to deploy your Spark cluster on either AWS or IBM Cloud
Elective 5: Your Choice	Use your skills to tackle any other project of your choice

T + 1 - e
• Introduction
Natural Language Processing
• Transformers, what can they do?
How do Transformers work?
• Encoder models
Decoder models
Sequence-to-sequence models Bigs and limitations
Bias and limitations
• Introduction
Behind the pipeline
• Models
• Tokenizers
Handling multiple sequences
• Introduction
Processing the data
Fine-tuning a model with the Trainer API or Keras
• A full training
The Hugging Fore Hub
The Hugging Face Hub Hair a protein and models
Using pretrained models Charing pretrained models
Sharing pretrained models Ruilding a model pard
Building a model card
• Introduction
What if my dataset isn't on the Hub?
Time to slice and dice
• Big data?
Creating your own dataset
Semantic search with FAISS
Learn the difference between Regression and Classification
Train a Linear Regression model to predict values
Learn to predict states using Logistic Regression
• Learn the definition of a perceptron as a building block for neural networks, and
the perceptron algorithm for classification.
Train Decision Trees to predict states
Use Entropy to build decision trees, recursively
Learn Bayes' rule, and apply it to predict cases of spam messages using the
Naive Bayes algorithm.
Train models using Bayesian Learning
Complete an exercise that uses Bayesian Learning for
Complete an exercise that uses Bayesian Learning for natural language processing

Ensemble of Learners	 Build data visualizations for quantitative and categorical data. Create pie, bar, line, scatter, histogram, and boxplot charts. Build professional presentations.
Evaluation Metrics	 Learn about different metrics to measure model success. Calculate accuracy, precision, and recall to measure the performance of your models.
Training and Tuning Models	Train and test models with Scikit-learn. Choose the best model using evaluation techniques like cross-validation and grid search.
Neural Networks	
Introduction to Neural Networks	Learn the foundations of deep learning and neural networks. Implement gradient descent and backpropagation in Python
Implementing Gradient Descent	Implement gradient descent using NumPy matrix multiplication.
Training Neural Networks	 Learn several techniques to effectively train a neural network. Prevent overfitting of training data and learn best practices for minimizing the error of a network.
Deep Learning with PyTorch	Learn how to use PyTorch for building deep learning models.
Unsupervised Learning	
Clustering	Learn the basics of clustering data Cluster data with the K-means algorithm
Hierarchical and Density-Based Clustering	 Cluster data with Single Linkage Clustering. Cluster data with DBSCAN, a clustering method that captures the insight that clusters are dense group of points
Gaussian Mixture Models	Cluster data with Gaussian Mixture Models Optimize Gaussian Mixture Models with and Expectation Maximization
Dimensionality Reduction	Reduce the dimensionality of the data using Principal Component Analysis and Independent Component Analysis

Capstone Projects

The participants will complete a capstone project featuring a science challenge for each stream (usually 1 or 2 weeks)

The capstone projects will be presented by the end of the stream completion period to the peers and to selected members of the pillars.

The potential capstone projects will be advertise during the stream. Participants are also encourage to propose their own ones.