# Course Outline for

Data Sciences/Optimization

Cloud Data Services – AWS, GCP and Azure with governance

Visualization and Charting

KG and Vector DBs

RDBMS and NoSQL

# Contents

# Stream 1

## Cloud Data Services for AWS, GCP, and Azure with Governance and Catalog

### Course Description:

This course provides a comprehensive overview of data services offered by the three major cloud providers: AWS, Google Cloud Platform (GCP), and Microsoft Azure. It covers essential data services, data governance, and data cataloging, emphasizing practical applications and hands-on experience. Students will learn to leverage cloud-based data services for storage, processing, and analysis while ensuring data governance and effective cataloging.

### 1: Introduction to Cloud Data Services

*Overview of Cloud Computing*
- o Cloud Computing Models (IaaS, PaaS, SaaS)
- o Benefits of Cloud Data Services

*Introduction to AWS, GCP, and Azure*
- o Key Differences and Similarities
- o Use Cases and Industry Adoption

*Setting Up Cloud Accounts*
- o Creating AWS, GCP, and Azure Accounts
- o Navigating Cloud Consoles

### 2: Data Storage Services

*AWS Data Storage Services*
- o Amazon S3 (Simple Storage Service)
- o Amazon EBS (Elastic Block Store)
- o Amazon Glacier

*GCP Data Storage Services*
- o Google Cloud Storage
- o Google Persistent Disk
- o Google Cloud Filestore

*Azure Data Storage Services*
- o Azure Blob Storage
- o Azure Disk Storage
- o Azure File Storage

*Hands-On Labs*

- o Setting Up and Managing Storage Services

## 3: Data Processing and Analytics

*AWS Data Processing and Analytics*

- o Amazon RDS (Relational Database Service)

- o Amazon Redshift

- o Amazon EMR (Elastic MapReduce)

*GCP Data Processing and Analytics*

- o Google BigQuery

- o Google Cloud SQL

- o Google Dataproc

*Azure Data Processing and Analytics*

- o Azure SQL Database

- o Azure Synapse Analytics

- o Azure HDInsight

*Hands-On Labs*

- o Deploying and Managing Data Processing Services

## 4: Big Data and Machine Learning

*AWS Big Data and ML Services*

- o AWS Glue

- o Amazon SageMaker

- o Amazon Kinesis

*GCP Big Data and ML Services*

- o Google Dataflow

- o Google AI Platform

- o Google Pub/Sub

*Azure Big Data and ML Services*

- o Azure Data Factory

- o Azure Machine Learning

- o Azure Stream Analytics

*Hands-On Labs*

- o Implementing Big Data and ML Solutions

## 5: Data Integration and ETL

*AWS Data Integration Services*

- o AWS Glue

- o AWS Data Pipeline

*GCP Data Integration Services*
- o Google Cloud Dataflow

- o Google Cloud Dataprep

*Azure Data Integration Services*
- o Azure Data Factory

*ETL Concepts and Tools*
- o Extract, Transform, Load (ETL) Processes

- o Comparing ETL Tools

*Hands-On Labs*
- o Creating ETL Pipelines

## 6: Data Governance Fundamentals

*Introduction to Data Governance*
- o Importance of Data Governance

- o Key Components of Data Governance

*AWS Data Governance*
- o AWS Lake Formation

- o AWS IAM (Identity and Access Management)

- o AWS Audit Manager

*GCP Data Governance*
- o Google Cloud Data Catalog

- o Google Cloud IAM

- o Google Cloud Audit Logs

*Azure Data Governance*
- o Azure Purview

- o Azure Active Directory

- o Azure Policy and Blueprints

## 7: Data Security and Compliance

*Security Best Practices*
- o Encryption and Key Management

- o Identity and Access Management

*AWS Security and Compliance*
- o AWS KMS (Key Management Service)

- o AWS Shield and WAF

*GCP Security and Compliance*
- o Google Cloud KMS

- o Google Cloud Armor

*Azure Security and Compliance*
- o Azure Key Vault

- o Azure Security Center

*Hands-On Labs*
- o Implementing Security Measures

## 8: Data Cataloguing and Metadata Management

*Introduction to Data Cataloging*
- o Importance of Data Cataloging

- o Key Features of Data Catalogs

*AWS Data Cataloging*
- o AWS Glue Data Catalog

*GCP Data Cataloging*
- o Google Cloud Data Catalog

*Azure Data Cataloging*
- o Azure Data Catalog (Deprecated)

- o Azure Purview

*Hands-On Labs*
- o Creating and Managing Data Catalogs

## 9: Advanced Data Governance and Policy Management

*Data Lineage and Ancestry*
- o Tracking Data Movement and Transformation

*Policy Management*
- o Defining and Enforcing Data Policies

*Data Quality Management*
- o Ensuring Data Accuracy and Consistency

*Case Studies and Best Practices*
- o Implementing Governance Frameworks

## 10: Data Lake and Data Warehouse Solutions

*Building Data Lakes*
- o Concepts and Architectures

- o AWS Lake Formation

- o Google Cloud Storage

- o Azure Data Lake Storage

*Data Warehousing*
- o Data Warehouse Concepts

- o Amazon Redshift

- o Google BigQuery

- o Azure Synapse Analytics

*Hands-On Labs*
- o Setting Up Data Lakes and Warehouses

## 11: Monitoring and Optimization

*Monitoring Data Services*
- o AWS CloudWatch

- o Google Cloud Monitoring

- o Azure Monitor

*Performance Optimization*
- o Scaling and Tuning Data Services

- o Cost Management and Optimization

*Hands-On Labs*
- o Implementing Monitoring Solutions

## 12: Capstone Project and Case Studies

*Capstone Project*
- o Designing a Comprehensive Cloud Data Solution

- o Implementing Data Governance and Cataloging

- o Ensuring Security and Compliance

- o Optimizing Performance and Cost

*Case Studies*
- o Real-World Applications of Cloud Data Services

- o Success Stories and Lessons Learned

*Final Presentations*
- o Presenting Capstone Projects

- o Peer Reviews and Feedback

## Software and Tools:

- o AWS Management Console

- o Google Cloud Console

- o Azure Portal

- o Data Integration and ETL Tools

- o Data Catalog and Governance Tools

*Assessment:*

- **Assignments:**

  - o ly assignments on cloud data services and governance

- **Projects:**

  - o Midterm project on a specific cloud platform

  - o Capstone project integrating all three platforms

- **Exams:**

  - o Midterm and final exams testing theoretical and practical knowledge

**Prerequisites:**

- Basic knowledge of cloud computing concepts

- Familiarity with data management principles

- Understanding of basic programming concepts (preferably in Python)

# Real-Time Data Processing

## Introduction to Real-Time Data Processing

- Overview of Real-Time Data Processing

- Use Cases and Applications

- Key Concepts: Stream Processing vs. Batch Processing

## Apache Kafka

- Introduction to Kafka

  - o Architecture and Components

  - o Producers, Consumers, and Brokers

- Setting Up Kafka

  - o Installation and Configuration

- Hands-On Lab: Building a Simple Kafka Producer and Consumer

## Real-Time Data Processing with Cloud Services

*AWS Kinesis*

- Overview of AWS Kinesis

  - o Kinesis Data Streams

  - o Kinesis Data Firehose

  - o Kinesis Data Analytics

- Setting Up Kinesis

    o Creating Data Streams

- Hands-On Lab: Stream Data with Kinesis

## GCP Pub/Sub and Azure Event Hub

- Introduction to GCP Pub/Sub

    o Architecture and Components

    o Setting Up Pub/Sub

    o Hands-On Lab: Publish and Subscribe to Messages

- Introduction to Azure Event Hub

    o Architecture and Components

    o Setting Up Event Hub

    o Hands-On Lab: Stream Data with Event Hub

# Big Data Processing and Data Governance

## Introduction to Spark and Databricks

- Overview of Apache Spark

    o Spark Architecture and Components

    o RDDs and DataFrames

- Introduction to Databricks

    o Databricks Environment

    o Integration with Spark

- Hands-On Lab: Basic Data Processing with Spark and Databricks

## Data Governance and Classification

- Basic Concepts of Data Classification

    o Types of Data: PII, Sensitive Data, etc.

    o Importance of Data Classification

- Introduction to Data Annotation

    o Tools and Techniques

- Using Atlas Tool for Data Governance

    o Overview and Features

    o Hands-On Lab: Data Annotation and Classification with Atlas Tool

# Data Warehousing and Data Lakes

## Introduction to Data Warehousing and Data Lakes

- Overview of Data Warehousing and Data Lakes

    o Key Concepts and Differences

- Introduction to Azure Synapse

    o Features and Architecture

    o Setting Up and Using Synapse

- Hands-On Lab: Basic Operations with Azure Synapse

## GCP Spanner and Alloy DB

- Introduction to GCP Spanner

    o Features and Architecture

    o Setting Up and Using Spanner

- Introduction to Alloy DB

    o Features and Architecture

    o Setting Up and Using Alloy DB

- Hands-On Lab: Basic Operations with GCP Spanner and Alloy DB

# Visualization and Charting with Data

## Course Description:

This course provides a comprehensive overview of data visualization and charting techniques using various tools and programming languages. Students will learn to create effective visualizations that convey data insights clearly and compellingly. The course covers foundational principles, advanced techniques, and best practices in data visualization.

## Course Outline:

### 1: Introduction to Data Visualization

- **Overview of Data Visualization**
    - Importance of Data Visualization
    - History and Evolution

- **Principles of Effective Visualization**
    - Visual Perception
    - Data-Ink Ratio
    - Choosing the Right Chart
    - Introduction to Co-pilot in Power BI
    - Overview and Capabilities
    - Setting Up Co-pilot in Power BI
    - Creating Dashboards with Co-pilot
    - Automated Insights and Suggestions
    - Customizing Visualizations with Co-pilot Assistance
    - Hands-On Lab: Building an Interactive Dashboard with Co-pilot
    - Importing Data
    - Creating Visuals and Charts
    - Applying Co-pilot Recommendations

- **Tools and Libraries Overview**
    - Overview of GCP Looker
    - Key Features and Advantages over Tableau
    - Integration with GCP Ecosystem
    - Setting Up GCP Looker
    - Account Setup and Configuration

- o Connecting Data Sources
- o Creating Visualizations in GCP Looker
- o Building Dashboards and Reports
- o Advanced Visuals and Interactivity
- o Python (Matplotlib, Seaborn)
- o JavaScript (D3.js)

## 2: Basic Chart Types

- **Bar Charts**
  - o Vertical and Horizontal Bars
  - o Grouped and Stacked Bars
- **Line Charts**
  - o Single and Multiple Lines
  - o Smoothing Techniques
- **Pie and Donut Charts**
  - o When to Use and When to Avoid
  - o Alternatives to Pie Charts
- **Hands-On Labs**
  - o Creating Basic Charts with Different Tools

## 3: Advanced Chart Types

- **Scatter Plots**
  - o Adding Trend Lines
  - o Bubble Charts
- **Histograms and Density Plots**
  - o Distribution Analysis
  - o Kernel Density Estimation
- **Box Plots and Violin Plots**
  - o Comparative Analysis
- **Hands-On Labs**
  - o Creating Advanced Charts with Python and R

## 4: Interactive Visualizations

- **Introduction to Interactivity**
  - o Importance and Applications

- **Interactive Tools and Libraries**

    o Plotly, Bokeh

    o Dash, Shiny

- **Building Interactive Dashboards**

    o Linking Multiple Visuals

    o Adding Filters and Controls

- **Hands-On Labs**

    o Creating Interactive Visualizations with Plotly and Dash

## 5: Data Visualization with D3.js

- **Introduction to D3.js**

    o Understanding the Basics

    o Selections and Data Binding

- **Creating Basic Visualizations**

    o Bar Charts, Line Charts, Scatter Plots

- **Advanced D3.js Techniques**

    o Transitions and Animations

    o Custom Visualizations

- **Hands-On Labs**

    o Building Visualizations with D3.js

## 6: Geospatial Visualization

- **Introduction to Geospatial Data**

    o Types and Sources of Geospatial Data

- **Mapping Tools and Libraries**

    o Leaflet, Folium

    o Mapbox, Google Maps API

- **Creating Maps and Geospatial Charts**

    o Choropleth Maps

    o Heatmaps and Marker Clusters

- **Hands-On Labs**

    o Visualizing Geospatial Data with Python and JavaScript

## 7: Visualization Best Practices

- **Design Principles**

- o Color Theory and Usage

- o Typography and Layout

- **Avoiding Common Pitfalls**

  - o Misleading Visualizations

  - o Overcomplicating Charts

- **Storytelling with Data**

  - o Crafting a Narrative

  - o Annotating Visualizations

- **Hands-On Labs**

  - o Redesigning Poor Visualizations

## 8: Data Visualization in Business Intelligence

- **Introduction to Business Intelligence (BI) Tools**

  - o Tableau, Power BI

- **Connecting to Data Sources**

  - o Databases, APIs, Spreadsheets

- **Creating Dashboards and Reports**

  - o Best Practices for BI Dashboards

  - o Real-Time Data Integration

- **Hands-On Labs**

  - o Building BI Dashboards with Tableau and Power BI

## 9: Advanced Visualization Techniques

- **Time Series Visualization**

  - o Line Charts, Area Charts, and Slope Graphs

- **Network Visualization**

  - o Node-Link Diagrams, Sankey Diagrams

- **Multidimensional Data Visualization**

  - o Parallel Coordinates, Radar Charts

- **Hands-On Labs**

  - o Implementing Advanced Techniques with Python and R

## 10: Capstone Project and Case Studies

- **Capstone Project**

  - o Selecting a Dataset

- o Designing and Implementing Visualizations

- o Presenting Findings and Insights

- **Case Studies**

  - o Analysis of Real-World Visualization Projects

  - o Success Stories and Lessons Learned

- **Final Presentations**

  - o Presenting Capstone Projects

  - o Peer Reviews and Feedback

- **Software and Tools:**

  - o Tableau, Power BI

  - o Python (Matplotlib, Seaborn, Plotly, Dash)

  - o JavaScript (D3.js)

  - o R (ggplot2, Shiny)

- **Graph Visualization using D3.js**

- **Introduction to Graph Visualization**

- Overview of Graphs and Networks

- Nodes and Edges

- Types of Graphs (Directed, Undirected, Weighted, etc.)

- Importance of Graph Visualization

- Use Cases and Applications (Knowledge Graphs, Social Networks, etc.)

- Introduction to D3.js

- Features and Capabilities

- Setting Up D3.js

- **Creating Graph Visualizations with D3.js**

- Basic Concepts in D3.js

- Selections and Data Binding

- SVG Elements (Circles, Lines, etc.)

- Building a Simple Graph

- Representing Nodes and Edges

- Adding Labels and Tooltips

- Hands-On Lab: Creating a Simple Graph Visualization

- Loading Data

- Rendering Nodes and Edges

- Adding Interactivity

- **Advanced Graph Visualization and Introduction to Graph Databases**

- **Advanced Graph Visualization Techniques**

- Force-Directed Layouts

- Understanding Force Simulations

- Configuring Forces (Link, Charge, Center, etc.)

- Enhancing Visualizations

- Styling and Animations

- Handling Large Graphs

- Hands-On Lab: Creating an Interactive Knowledge Graph

- Implementing Force-Directed Layout

- Adding Interactivity and Animations

- **Introduction to Azure Cosmos DB and AWS Neptune**

- Overview of Graph Databases

- What are Graph Databases?

- Use Cases and Applications

- Introduction to Azure Cosmos DB

- Features and Capabilities

- Data Model and API (Gremlin API)

- Introduction to AWS Neptune

- Features and Capabilities

- Data Model and API (Gremlin and SPARQL)

- Hands-On Lab: Setting Up and Querying Graph Databases

- Setting Up Azure Cosmos DB

- Setting Up AWS Neptune

- Basic CRUD Operations

- Writing and Executing Graph Queries


## Assessment:

- **Assignments:**
    - ly assignments on creating various types of visualizations

## Projects:

- o Midterm project on interactive visualizations

- o Capstone project integrating multiple visualization techniques

# Knowledge Graphs and Vector Databases

## Course Description:

This course provides a comprehensive understanding of Knowledge Graphs (KG) and Vector Databases (Vector DBs). Students will learn the principles, techniques, and applications of these technologies in managing and querying complex, interconnected data. The course covers foundational concepts, practical implementation, and advanced topics, emphasizing hands-on experience with relevant tools and platforms.

## Course Outline:

### 1: Introduction to Knowledge Graphs

- **Overview of Knowledge Graphs**
  - Definition and History
  - Importance and Applications
- **Fundamental Concepts**
  - Nodes, Edges, and Properties
  - Ontologies and Taxonomies
  - RDF (Resource Description Framework) and SPARQL
- **Tools and Platforms**
  - Neo4j, Apache Jena, GraphDB
- **Hands-On Labs**
  - Creating a Simple Knowledge Graph

### 2: Data Modeling with Knowledge Graphs

- **Data Modeling Principles**
  - Entity-Relationship Modeling
  - Schema Design for Knowledge Graphs
- **Ontology Development**
  - Creating and Using Ontologies
  - OWL (Web Ontology Language)
- **Case Studies**
  - Real-World Examples of Knowledge Graphs
- **Hands-On Labs**
  - Designing and Implementing a Knowledge Graph Schema

## 3: Querying Knowledge Graphs

- **SPARQL Fundamentals**

    o Basic SPARQL Queries

    o Filtering and Aggregation

- **Advanced SPARQL Techniques**

    o Joins and Subqueries

    o SPARQL Extensions (SPIN, SHACL)

- **Hands-On Labs**

    o Writing and Executing SPARQL Queries

## 4: Knowledge Graph Construction and Integration

- **Data Ingestion Techniques**

    o ETL (Extract, Transform, Load) Processes

    o Data Integration from Various Sources

- **Knowledge Graph Enrichment**

    o Entity Linking and Resolution

    o Semantic Annotation

- **Hands-On Labs**

    o Building and Enriching a Knowledge Graph

## 5: Introduction to Vector Databases

- **Overview of Vector Databases**

    o Definition and Applications

    o Key Differences from Traditional Databases

- **Fundamental Concepts**

    o Vector Representation of Data

    o Embedding Spaces and Similarity Search

- **Tools and Platforms**

    o FAISS, Milvus, Pinecone

- **Hands-On Labs**

    o Setting Up and Querying a Vector Database

## 6: Data Modelling with Vector Databases

- **Vector Representation Techniques**

    o Word Embeddings (Word2Vec, GloVe)

- o Sentence Embeddings (BERT, GPT)

- **Indexing and Retrieval**

    - o Approximate Nearest Neighbor (ANN) Search

    - o Index Structures (LSH, HNSW)

- **Case Studies**

    - o Real-World Examples of Vector Databases

- **Hands-On Labs**

    - o Creating and Querying Embeddings in a Vector Database

# 7: Advanced Topics in Knowledge Graphs

- **Graph Algorithms**

    - o Pathfinding and Shortest Path

    - o Centrality Measures

    - o Community Detection

- **Scalability and Performance**

    - o Distributed Knowledge Graphs

    - o Optimization Techniques

- **Hands-On Labs**

    - o Implementing and Analyzing Graph Algorithms

# 8: Advanced Topics in Vector Databases

- **Dimensionality Reduction**

    - o PCA, t-SNE, UMAP

- **Hybrid Search Techniques**

    - o Combining Vector and Traditional Queries

- **Scalability and Performance**

    - o Distributed Vector Databases

    - o Optimization Techniques

- **Hands-On Labs**

    - o Applying Dimensionality Reduction and Hybrid Search

# 9: Integration of Knowledge Graphs and Vector Databases

- **Complementary Use Cases**

    - o Enhancing Knowledge Graphs with Embeddings

    - o Leveraging Vector Databases for Semantic Search

- **Integration Techniques**

    o Data Pipelines and ETL for Combined Use

    o Querying Across Knowledge Graphs and Vector Databases

- **Hands-On Labs**

    o Integrating and Querying Combined Systems

## 10: Capstone Project and Case Studies

- **Capstone Project**

    o Selecting a Real-World Problem

    o Designing and Implementing a Solution Using Knowledge Graphs and Vector Databases

    o Presenting Findings and Insights

- **Case Studies**

    o Analysis of Successful Implementations

    o Challenges and Best Practices

- **Final Presentations**

    o Presenting Capstone Projects

    o Peer Reviews and Feedback

- **Software and Tools:**

    o Neo4j, Apache Jena, GraphDB

    o FAISS, Milvus, Pinecone

    o Python Libraries: NetworkX, Gensim, scikit-learn

## Assessment:

- **Assignments:**

    o ly assignments on various aspects of knowledge graphs and vector databases

- **Projects:**

    o Midterm project on constructing a knowledge graph

    o Capstone project integrating knowledge graphs and vector databases

## Prerequisites:

- Basic knowledge of data structures and algorithms

- Familiarity with database concepts

- Understanding of basic programming concepts (preferably in Python)

# Relational Database Management Systems (RDBMS) and NoSQL Databases

## Course Description:

This course provides a comprehensive overview of Relational Database Management Systems (RDBMS) and NoSQL databases. Students will learn the fundamentals of database design, implementation, and management for both relational and non-relational databases. The course covers SQL and NoSQL query languages, data modeling, and the advantages and use cases of each database type.

## Course Outline:

### 1: Introduction to Databases

- **Overview of Databases**
    - History and Evolution of Databases
    - Types of Databases

- **Introduction to RDBMS**
    - Key Concepts: Tables, Rows, Columns, Primary Keys, Foreign Keys
    - ACID Properties

- **Introduction to NoSQL**
    - Key Concepts: Collections, Documents, Key-Value Pairs
    - BASE Properties

### 2: Database Design and Data Modeling (RDBMS)

- **Entity-Relationship (ER) Modeling**
    - Entities, Attributes, and Relationships
    - ER Diagrams

- **Normalization**
    - Normal Forms (1NF, 2NF, 3NF, BCNF)
    - De-normalization

- **SQL Data Definition Language (DDL)**
    - Creating and Modifying Tables
    - Constraints and Indexes

### 3: SQL Basics

- **SQL Data Manipulation Language (DML)**
    - SELECT, INSERT, UPDATE, DELETE

- **SQL Query Syntax**

  - Basic Queries

  - Filtering and Sorting Data

- **Joins and Subqueries**

  - INNER JOIN, LEFT JOIN, RIGHT JOIN, FULL JOIN

  - Nested Queries

## 4: Advanced SQL

- **Advanced Query Techniques**

  - Aggregate Functions (SUM, AVG, COUNT, etc.)

  - GROUP BY and HAVING Clauses

- **Stored Procedures and Functions**

  - Creating and Using Stored Procedures

  - Creating and Using Functions

- **Transactions and Concurrency Control**

  - BEGIN, COMMIT, ROLLBACK

  - Isolation Levels

## 5: Introduction to NoSQL Databases

- **Types of NoSQL Databases**

  - Document Stores (e.g., MongoDB)

  - Key-Value Stores (e.g., Redis)

  - Column Stores (e.g., Cassandra)

  - Graph Databases (e.g., Neo4j)

- **NoSQL vs. RDBMS**

  - Scalability and Performance

  - Use Cases and Trade-offs

## 6: Document-Oriented Databases (MongoDB)

- **MongoDB Basics**

  - Collections and Documents

  - BSON Format

- **CRUD Operations**

  - Inserting, Updating, Deleting Documents

  - Querying Documents

- **Indexing and Aggregation**

  - Creating Indexes

  - Aggregation Framework

# 7: Key-Value Stores (Redis)

- **Redis Basics**

  - Key-Value Data Model

  - Data Types in Redis (Strings, Lists, Sets, Hashes, etc.)

- **Redis Commands**

  - Basic CRUD Operations

  - Transactions and Pipelining

- **Use Cases**

  - Caching, Session Management, Real-Time Analytics

# 8: Column-Oriented Databases (Apache Cassandra)

- **Cassandra Basics**

  - Data Model (Keyspace, Column Families, Rows)

  - Consistency Levels

- **Cassandra Query Language (CQL)**

  - Basic CRUD Operations

  - Advanced CQL Queries

- **Data Replication and Consistency**

  - Replication Strategies

  - Handling Consistency

# 9: Graph Databases (Neo4j)

- **Graph Database Basics**

  - Nodes, Relationships, Properties

  - Graph Data Model

- **Cypher Query Language**

  - Basic CRUD Operations

  - Traversing Graphs

- **Use Cases**

  - Social Networks, Recommendation Systems, Fraud Detection

## 10: Data Management and Security

- **Backup and Recovery**

    o Strategies for Backup and Recovery

    o Tools and Techniques

- **Security**

    o Authentication and Authorization

    o Encryption and Data Masking

- **Performance Tuning**

    o Indexing and Query Optimization

    o Monitoring and Diagnostics

## 11: Database Integration and Migration

- **Integrating RDBMS and NoSQL**

    o Hybrid Database Systems

    o Data Sync and ETL Processes

- **Database Migration**

    o Migrating from RDBMS to NoSQL

    o Tools and Best Practices

- **Case Studies**

    o Real-World Examples of Database Integration and Migration

## 12: Capstone Project and Case Studies

- **Capstone Project**

    o Designing and Implementing a Database Solution

    o Integrating RDBMS and NoSQL Components

    o Presenting Project Findings

- **Case Studies**

    o Analysis of Industry Use Cases

    o Lessons Learned from Real-World Implementations

- **Final Presentations**

    o Presenting Capstone Projects

    o Peer Reviews and Feedback

- **Software and Tools:**

    o MySQL, PostgreSQL

- o MongoDB, Redis, Apache Cassandra, Neo4j

- o SQL Management Tools (e.g., MySQL Workbench, pgAdmin)

- o NoSQL Management Tools (e.g., MongoDB Compass, RedisInsight)

## Assessment:

- **Assignments:**

  - o ly assignments on database design, SQL queries, and NoSQL operations

- **Projects:**

  - o Midterm project on RDBMS design and implementation

  - o Capstone project integrating RDBMS and NoSQL solutions

  **Prerequisites:**

- Basic knowledge of programming concepts

- Familiarity with data structures and algorithms

- Understanding of basic computer science principles

# Hadoop

## Day 1: Introduction to Hadoop and HDFS

- Overview of Big Data

- History and Evolution of Hadoop

- Hadoop Ecosystem Components

## Hadoop Distributed File System (HDFS)

## HDFS Architecture

- NameNode and DataNode

- Block Storage and Replication

## HDFS Commands

- Basic File Operations (put, get, list, etc.)

- Hands-On Lab: Setting Up a Hadoop Environment

  - Installing Hadoop

  - Configuring HDFS

  - Running Basic HDFS Commands

## Introduction to MapReduce

- MapReduce Framework

  - Map and Reduce Functions

  - Data Flow in MapReduce

- MapReduce Example

  - Word Count Program

## YARN (Yet Another Resource Negotiator)

- YARN Architecture

  - ResourceManager and NodeManager

  - ApplicationMaster and Containers

- Resource Allocation and Job Scheduling

## Hands-On Lab: Writing and Running MapReduce Jobs

- Setting Up a Sample Data Set

- Writing a Simple MapReduce Program

- Submitting and Monitoring Jobs

# Introduction to Spark and PySpark

# Introduction to Apache Spark

## Spark Overview

- o Differences Between Hadoop and Spark

- o Spark Ecosystem Components

## Spark Architecture

- o Driver and Executors

- o RDDs (Resilient Distributed Datasets)

- o DAG (Directed Acyclic Graph)

# Getting Started with PySpark

## Setting Up PySpark

- o Installing and Configuring PySpark

- o PySpark Shell and Notebooks

## Basic PySpark Operations

- o Creating RDDs

- o Transformations and Actions

## DataFrame API in PySpark

- Introduction to DataFrames

  - o Creating DataFrames

  - o Performing Operations on DataFrames

- SQL Queries with DataFrames

  - o Using Spark SQL

- Hands-On Lab: Working with DataFrames in PySpark

  - o Loading Data into DataFrames

  - o Applying Transformations and Actions

## Spark Streaming

- Introduction to Spark Streaming

  - o Streaming Architecture

  - o DStreams (Discretized Streams)

- Processing Real-Time Data

- o Window Operations

- o Integrating with Data Sources (Kafka, HDFS, etc.)

## Machine Learning with MLlib

- Introduction to Spark MLlib

  - o MLlib vs. Spark ML

  - o Pipelines and Estimators

- Example Use Case: Predictive Modeling

  - o Loading and Preparing Data

  - o Training and Evaluating Models

## Hands-On Lab: End-to-End Data Processing with PySpark

- Loading Data from HDFS

- Transforming Data with DataFrame API

- Running a Machine Learning Pipeline

- Streaming Data Processing


## Assessment:

- **Hands-On Labs:**

  - o Practical exercises on HDFS, MapReduce, and PySpark

- **Quizzes:**

  - o Short quizzes to test understanding of key concepts

- **Capstone Project:**

  - o An end-to-end data processing project using Hadoop and Spark

## Prerequisites:

- Basic knowledge of programming (preferably Python)

- Familiarity with Linux command line

- Understanding of basic data processing concepts

# Stream 2

## Data Science and Python

### Course Outline:

This course aims to provide an in-depth understanding of data science principles and optimization techniques using Python. Students will learn to manipulate and analyze data, build predictive models, and implement optimization algorithms. The course combines theoretical concepts with practical applications, ensuring that students can apply their knowledge to real-world problems.

### 1: Introduction to Data Science and Python

### Introduction to Data Science

- o What is Data Science?
- o Applications of Data Science
- o The Data Science Process

### Python Basics

- o Introduction to Python
- o Data Types, Variables, and Operators
- o Control Structures (Loops and Conditionals)
- o Functions and Modules

### Python Libraries for Data Science

- o NumPy, Pandas, Matplotlib, and Seaborn

### 2: Data Manipulation and Analysis with Pandas

*Introduction to Pandas*

- o Series and DataFrames
- o Importing and Exporting Data

*Data Cleaning and Preparation*

- o Handling Missing Data
- o Data Transformation
- o Data Aggregation and Grouping

*Exploratory Data Analysis (EDA)*

- o Descriptive Statistics
- o Data Visualization Techniques

## 3: Data Visualization

*Introduction to Data Visualization*

- o Importance of Data Visualization

- o Types of Visualizations

*Matplotlib and Seaborn*

- o Basic Plotting with Matplotlib

- o Advanced Visualization with Seaborn

- o Customizing Plots

*Interactive Visualization with Plotly*

- o Introduction to Plotly

- o Creating Interactive Dashboards

## 4: Probability and Statistics

- **Introduction to Probability**

  - o Basic Probability Concepts

  - o Probability Distributions

- **Statistical Inference**

  - o Descriptive vs. Inferential Statistics

  - o Hypothesis Testing

  - o Confidence Intervals

## 5: Machine Learning Basics

- **Introduction to Machine Learning**

  - o Types of Machine Learning

  - o Supervised vs. Unsupervised Learning

- **Supervised Learning**

  - o Linear Regression

  - o Logistic Regression

  - o Evaluation Metrics (RMSE, Accuracy, Precision, Recall)

## 6: Advanced Machine Learning

- **Unsupervised Learning**

  - o Clustering (K-Means, Hierarchical)

  - o Dimensionality Reduction (PCA)

- **Advanced Algorithms**

  - o Decision Trees and Random Forests

- o Support Vector Machines

- o Ensemble Methods

# 7: Optimization Fundamentals

- **Introduction to Optimization**

  - o What is Optimization?

  - o Applications of Optimization

- **Mathematical Foundations**

  - o Linear Algebra and Calculus for Optimization

  - o Convex and Non-Convex Functions

- **Linear Programming**

  - o Formulating Linear Programs

  - o Simplex Method

  - o Duality in Linear Programming

# 8: Non-Linear and Integer Programming

- **Non-Linear Optimization**

  - o Gradient Descent

  - o Constrained Optimization (Lagrange Multipliers)

- **Integer Programming**

  - o Formulating Integer Programs

  - o Branch and Bound Method

# 9: Heuristic and Metaheuristic Methods

- **Introduction to Heuristics**

  - o Basic Concepts

  - o Greedy Algorithms

- **Metaheuristic Algorithms**

  - o Genetic Algorithms

  - o Simulated Annealing

  - o Particle Swarm Optimization

# 10: Optimization in Machine Learning

- **Optimization in Model Training**

  - o Loss Functions

  - o Gradient-Based Optimization (SGD, Adam)

- **Hyperparameter Tuning**
    - Grid Search
    - Random Search
    - Bayesian Optimization

## 11: Advanced Topics in Data Science and Optimization

- **Time Series Analysis**
    - Introduction to Time Series
    - ARIMA and Exponential Smoothing
- **Natural Language Processing (NLP)**
    - Text Preprocessing
    - Sentiment Analysis
- **Deep Learning**
    - Introduction to Neural Networks
    - Convolutional Neural Networks (CNNs)
    - Recurrent Neural Networks (RNNs)

# 12. Introduction to Probability and Statistics (30 minutes)

- Basic Concepts of Probability
- Overview of Statistical Inference

## Bayes Theorem

- Introduction to Bayes Theorem
    - Definition and Formula
    - Conditional Probability
- Applications of Bayes Theorem
    - Diagnostic Testing
    - Spam Filtering
    - Bayesian Inference
- Worked Examples
    - Calculating Posterior Probabilities
    - Examples in Real-World Contexts
- Hands-On Lab: Implementing Bayes Theorem
    - Practical Exercises with Data Sets
    - Using Python/R for Bayesian Calculations

## Markov Theorem

- Introduction to Markov Processes

    o Definition and Characteristics

    o Types of Markov Chains (Discrete and Continuous)

- Understanding the Markov Property

    o Memoryless Property

    o Transition Matrices

- Applications of Markov Theorem

    o Predictive Modeling

    o PageRank Algorithm

    o Stock Market Analysis

- Worked Examples

    o Constructing Transition Matrices

    o Solving for Steady-State Probabilities

- Hands-On Lab: Implementing Markov Chains

    o Practical Exercises with Transition Matrices

    o Using Python/R for Markov Chain Simulations

## . Advanced Topics in Bayesian and Markov Theorems

- Bayesian Networks

    o Structure and Inference

    o Applications in Machine Learning

- Hidden Markov Models (HMM)

    o Definition and Applications

    o Viterbi Algorithm

- Hands-On Lab: Advanced Implementations

    o Building Bayesian Networks

    o Implementing HMMs in Python/R

## Case Studies and Real-World Applications Case Study: Medical Diagnosis Using Bayesian Methods

- Case Study: Predictive Text Input Using Markov Chains

- Group Discussion: Practical Applications in Industry

- **Software and Tools:**

o Python (with libraries such as NumPy, SciPy, and PyMC3)

o R (with packages such as BayesianTools, markovchain)

o Jupyter Notebooks/RStudio

**Assessment:**

- **Hands-On Labs:**

  o Practical exercises on implementing Bayes Theorem and Markov Chains

  o Simulating Bayesian and Markov models using Python/R

- **Quizzes:**

  o Short quizzes to test understanding of key concepts

- **Final Project:**

  o An end-to-end project involving the application of Bayes Theorem and Markov Chains to a real-world data set

**Prerequisites:**

- Basic knowledge of probability and statistics

- Familiarity with Python or R programming

- Understanding of basic data analysis techniques

## 12: Project and Case Studies

- **Real-World Data Science Project**

  o Project Planning and Data Collection

  o Data Analysis and Model Building

  o Optimization and Results Presentation

- **Case Studies**

  o Case Study 1: Optimization in Supply Chain Management

  o Case Study 2: Predictive Maintenance using Machine Learning

  o Case Study 3: Optimizing Marketing Campaigns

- **Software and Tools:**

  o Python (Anaconda Distribution)

  o Jupyter Notebooks

  o Libraries: NumPy, Pandas, Matplotlib, Seaborn, Scikit-Learn, Plotly

**Prerequisites:**

- Basic knowledge of Python programming

- Understanding of basic mathematical concepts (algebra, calculus, probability)