



ISDM (INDEPENDENT SKILL DEVELOPMENT MISSION)



BIG DATA CONCEPTS: HADOOP, SPARK, NoSQL DATABASES

📌 CHAPTER 1: INTRODUCTION TO BIG DATA

1.1 What is Big Data?

Big Data refers to **large and complex datasets** that traditional databases cannot handle efficiently. These datasets require advanced storage, processing, and analytical techniques to extract valuable insights.

Big Data is characterized by the **5 Vs**:

- ✓ **Volume** – The enormous amount of data generated daily (e.g., social media posts, IoT sensor data).
- ✓ **Velocity** – The speed at which data is created and processed (e.g., real-time stock market data).
- ✓ **Variety** – Different types of data (structured, unstructured, semi-structured).
- ✓ **Veracity** – The trustworthiness and reliability of data.
- ✓ **Value** – The usefulness of data in decision-making.

1.2 Why is Big Data Important?

Big Data is transforming industries by enabling:

- ✓ **Personalized marketing and customer experience** (Amazon, Netflix recommendations).
- ✓ **Real-time fraud detection** in banking.
- ✓ **Improved healthcare diagnostics** with AI-based predictions.
- ✓ **Optimized supply chain and logistics management.**

📌 Example Use Case:

A retail company tracks customer purchase history, website behavior, and social media interactions. Using Big Data analytics, they can **offer personalized recommendations, predict shopping trends, and enhance customer engagement.**

💡 Conclusion:

Big Data is a **game-changer** that helps businesses **optimize operations, reduce costs, and improve decision-making.**

📌 CHAPTER 2: INTRODUCTION TO HADOOP

2.1 What is Hadoop?

Hadoop is an **open-source framework** for storing and processing Big Data using a distributed computing approach. It allows businesses to **handle massive datasets across multiple machines in parallel.**

2.2 Components of Hadoop

- ◆ **1. Hadoop Distributed File System (HDFS)**
 - A **scalable and fault-tolerant** storage system that splits large files into **blocks** and distributes them across multiple nodes.
 - Data is stored redundantly across nodes for reliability.
- ◆ **2. YARN (Yet Another Resource Negotiator)**

- Manages cluster resources and schedules tasks efficiently.
- Allocates CPU and memory to processing tasks dynamically.

◆ **3. MapReduce (Processing Framework)**

- A programming model for **processing large datasets in parallel**.
- **Map phase** filters and sorts data, while the **Reduce phase** aggregates results.

◆ **4. Hadoop Common**

- Contains libraries and utilities required by other Hadoop modules.

2.3 How Hadoop Works?

- ✓ **Step 1:** Data is divided into **blocks** and distributed across nodes in an HDFS cluster.
- ✓ **Step 2:** A MapReduce job is initiated to process data in **parallel** across nodes.
- ✓ **Step 3:** The final output is aggregated and stored back in HDFS.

📌 **Example Use Case:**

A telecom company processes **call logs from millions of users** daily using Hadoop to detect network anomalies and optimize bandwidth allocation.

💡 **Conclusion:**

Hadoop enables **scalable, distributed, and fault-tolerant** processing of Big Data, making it a key player in Big Data ecosystems.

📌 CHAPTER 3: INTRODUCTION TO APACHE SPARK

3.1 What is Apache Spark?

Apache Spark is a **fast and general-purpose Big Data framework** for real-time and batch data processing. It is **100x faster than Hadoop MapReduce** and supports various workloads like machine learning, graph processing, and streaming analytics.

3.2 Key Features of Apache Spark

- ✓ **In-Memory Processing** – Uses RAM instead of disk, making it extremely fast.
- ✓ **Supports Multiple Languages** – Works with Python, Scala, Java, and R.
- ✓ **Unified Analytics Platform** – Handles SQL queries, machine learning, and real-time streaming.
- ✓ **Fault Tolerance** – Automatically recovers from failures.

3.3 Apache Spark Components

- ◆ **1. Spark Core**
 - Manages memory, job execution, and fault recovery.
- ◆ **2. Spark SQL**
 - Supports querying structured data using SQL syntax.
- ◆ **3. Spark Streaming**
 - Processes real-time streaming data from sources like Twitter, Kafka, and IoT devices.
- ◆ **4. MLlib (Machine Learning Library)**
 - Provides built-in algorithms for classification, regression, clustering, and recommendation systems.

◆ 5. GraphX

- Handles large-scale graph processing (e.g., social network analysis).

📌 Example Use Case:

A financial institution uses Apache Spark to analyze **real-time transactions for fraud detection**, allowing instant alerts and prevention of fraudulent activities.

💡 Conclusion:

Apache Spark is **faster, scalable, and suitable for real-time applications**, making it ideal for handling modern Big Data workloads.

📌 CHAPTER 4: NO-SQL DATABASES FOR BIG DATA

4.1 What is NoSQL?

NoSQL (Not Only SQL) databases are **non-relational databases** designed for **scalable, high-performance data storage**. They handle **unstructured, semi-structured, and structured data** efficiently.

4.2 Types of NoSQL Databases

- ◆ **1. Key-Value Stores (e.g., Redis, DynamoDB)**
 - Store data as **key-value pairs**.
 - Ideal for caching and real-time analytics.
- ◆ **2. Document-Oriented Databases (e.g., MongoDB, CouchDB)**
 - Store data as **JSON or BSON documents**.
 - Used for **content management systems and product catalogs**.

◆ **3. Column-Family Stores (e.g., Cassandra, HBase)**

- Store data in **columns instead of rows** for fast read/write operations.
- Used in **distributed applications and time-series data analysis**.

◆ **4. Graph Databases (e.g., Neo4j, Amazon Neptune)**

- Store relationships between data as **nodes and edges**.
- Used in **social networks, fraud detection, and recommendation systems**.

📌 **Example Use Case:**

An e-commerce platform uses **MongoDB** to store customer profiles, purchase history, and product reviews efficiently.

💡 **Conclusion:**

NoSQL databases are **scalable, flexible, and suitable for handling unstructured data**, making them a key component of Big Data systems.

📌 **CHAPTER 5: BIG DATA PROCESSING WORKFLOW**

- ✓ **Step 1: Data Ingestion** – Collecting raw data from multiple sources (social media, IoT devices, transactional databases).
- ✓ **Step 2: Storage & Management** – Using HDFS, NoSQL, or cloud storage solutions.
- ✓ **Step 3: Data Processing** – Using Spark or MapReduce to transform and analyze data.
- ✓ **Step 4: Data Analysis & Visualization** – Using BI tools like Tableau, Power BI, or ML models for insights.

📌 **Example Use Case:**

A healthcare provider processes **electronic health records** from multiple hospitals to predict disease outbreaks.

💡 **Conclusion:**

A well-structured **Big Data workflow** ensures efficient data collection, storage, processing, and decision-making.

📌 **CHAPTER 6: FUTURE TRENDS IN BIG DATA**

- ✓ **Edge Computing** – Processing data closer to the source (IoT devices).
- ✓ **AI-Driven Big Data Analytics** – Using ML for automated decision-making.
- ✓ **Cloud-Based Big Data Solutions** – AWS, Google Cloud, and Azure Big Data services.
- ✓ **Real-Time Data Processing** – Streaming analytics for instant insights.

📌 **Example Use Case:**

Self-driving cars use **real-time Big Data processing** to detect obstacles and make driving decisions instantly.

💡 **Conclusion:**

Big Data is evolving rapidly, and organizations that **adopt modern technologies** will gain a competitive advantage in the data-driven world.

CLOUD PLATFORMS FOR DATA SCIENCE: AWS, GOOGLE CLOUD, AZURE

📌 CHAPTER 1: INTRODUCTION TO CLOUD COMPUTING FOR DATA SCIENCE

1.1 What is Cloud Computing?

Cloud computing is the **on-demand delivery of computing resources** (such as servers, storage, databases, networking, and analytics) over the internet. Instead of maintaining **physical hardware and software**, users can access cloud services on a **pay-as-you-go** basis.

◆ Why Cloud Computing for Data Science?

- ✓ **Scalability:** Scale computing power up or down based on demand.
- ✓ **Cost-Effective:** Pay only for the resources used, reducing upfront infrastructure costs.
- ✓ **Flexibility & Accessibility:** Access resources and datasets from anywhere.
- ✓ **Performance Optimization:** Cloud platforms provide high-performance GPUs, TPUs, and AI-optimized computing.

1.2 Key Benefits for Data Science

- ✓ **High-Speed Data Processing:** Enables big data analytics and real-time insights.
- ✓ **Collaborative Environment:** Supports team-based data science projects.
- ✓ **Machine Learning Integration:** Offers built-in AI/ML tools and pre-trained models.

✓ **Security & Compliance:** Ensures encrypted storage and compliance with industry regulations.

📌 **Example Use Case:**

A data science team working on **real-time fraud detection** can use cloud computing to process large transaction datasets without investing in costly infrastructure.

💡 **Conclusion:**

Cloud computing has revolutionized data science, enabling faster, more scalable, and more cost-efficient analytics solutions.

📌 CHAPTER 2: OVERVIEW OF CLOUD PLATFORMS FOR DATA SCIENCE

2.1 Leading Cloud Platforms

The three most widely used cloud platforms for Data Science are:

- ✓ **Amazon Web Services (AWS)** – Provides a full suite of analytics and AI tools.
- ✓ **Google Cloud Platform (GCP)** – Offers seamless integration with TensorFlow and AI-driven applications.
- ✓ **Microsoft Azure** – Focuses on enterprise-scale AI, machine learning, and big data analytics.

2.2 Cloud Computing Models for Data Science

Cloud Model	Description	Example
IaaS (Infrastructure as a Service)	Provides virtual machines, networking, and storage.	AWS EC2, Azure Virtual Machines

PaaS (Platform as a Service)	Provides a managed development environment for AI & ML applications.	Google AI Platform, AWS SageMaker
SaaS (Software as a Service)	Provides ready-to-use cloud-based applications.	Google BigQuery, Azure Machine Learning

📌 **Example Use Case:**

A startup working on AI-powered medical diagnosis **uses Google Cloud AI Platform** to train and deploy deep learning models for detecting diseases in medical images.

💡 **Conclusion:**

Each cloud platform offers unique advantages tailored to different business needs, from scalable computing to end-to-end ML pipelines.

📌 **CHAPTER 3: AMAZON WEB SERVICES (AWS) FOR DATA SCIENCE**

3.1 Overview of AWS for Data Science

Amazon Web Services (AWS) is the **most widely adopted cloud platform** for big data processing, AI, and machine learning. It offers **high-performance computing (HPC)** resources optimized for data science workflows.

3.2 Key AWS Services for Data Science

Service	Description
Amazon S3 (Simple Storage Service)	Stores structured and unstructured data efficiently.

Amazon SageMaker	Fully managed ML service for building, training, and deploying models.
AWS Lambda	Serverless computing for event-driven data processing.
Amazon Redshift	Data warehousing solution for big data analytics.
AWS Glue	ETL (Extract, Transform, Load) tool for processing raw data.
Amazon EMR	Hadoop-based big data processing engine.

3.3 Running a Machine Learning Model on AWS SageMaker

AWS SageMaker simplifies the ML workflow from data preprocessing to model deployment.

- ➡ **Steps to Run ML Model on AWS SageMaker:**
- 1 **Upload Data to S3:** Store datasets securely.
- 2 **Train ML Model:** Use built-in ML algorithms or custom TensorFlow/PyTorch models.
- 3 **Deploy the Model:** Use SageMaker endpoints to deploy models for inference.
- 4 **Monitor Performance:** Use AWS CloudWatch to analyze model efficiency.

💡 Conclusion:

AWS provides a **highly scalable and feature-rich** ecosystem for data scientists and AI researchers.

📌 CHAPTER 4: GOOGLE CLOUD PLATFORM (GCP) FOR DATA SCIENCE

4.1 Overview of Google Cloud for Data Science

Google Cloud Platform (GCP) is widely used for **AI/ML research, big data analytics, and deep learning** applications. It seamlessly integrates with **TensorFlow, Kubernetes, and Google AI tools**.

4.2 Key GCP Services for Data Science

Service	Description
Google BigQuery	Serverless data warehouse for analytics.
Google AI Platform	End-to-end ML development environment.
Google Cloud AutoML	No-code ML model training for non-experts.
Google Cloud Functions	Event-driven serverless functions.
TensorFlow Enterprise	Optimized TensorFlow for large-scale AI applications.

4.3 Deploying a ML Model Using Google Cloud AI Platform

📌 Steps to Deploy a Model:

- 1 **Store Dataset in BigQuery:** Load and preprocess data.
- 2 **Train Model in AI Platform:** Choose AutoML or custom deep learning models.
- 3 **Deploy as API Endpoint:** Use REST API for real-time predictions.
- 4 **Monitor with Cloud AI Metrics:** Analyze accuracy and model drift.

💡 Conclusion:

Google Cloud is **ideal for AI-first organizations**, offering deep learning frameworks and advanced big data tools.

📌 CHAPTER 5: MICROSOFT AZURE FOR DATA SCIENCE

5.1 Overview of Microsoft Azure for Data Science

Microsoft Azure provides enterprise-grade **AI, machine learning, and data processing** capabilities with a focus on business applications and automation.

5.2 Key Azure Services for Data Science

Service	Description
Azure Machine Learning	Fully managed ML model development and deployment.
Azure Synapse Analytics	Integrated analytics service for big data.
Azure Databricks	Apache Spark-based analytics engine.
Azure Cognitive Services	Pre-built AI models for vision, speech, and text analysis.
Azure Blob Storage	Scalable data storage for structured/unstructured data.

5.3 Building a Data Science Pipeline on Azure

📌 Steps to Set Up a Pipeline:

- 1 Store Data in Azure Blob Storage.
- 2 Use Azure Databricks for Data Processing.
- 3 Train Machine Learning Models in Azure ML.
- 4 Deploy AI Solutions using Azure Cognitive Services.

💡 Conclusion:

Azure is a **powerful choice for enterprises** looking for cloud-based **data science automation and AI-driven decision-making**.

📌 CHAPTER 6: COMPARISON OF AWS, GCP, AND AZURE FOR DATA SCIENCE

Feature	AWS	GCP	Azure
Best For	Scalability, AI/ML Integration	AI & Deep Learning	Enterprise AI Solutions
ML Platform	SageMaker	AI Platform	Azure ML
Big Data Services	Redshift, EMR	BigQuery	Synapse Analytics
AI Tools	AWS AI Services	TensorFlow AI	Cognitive Services
Cost-Effectiveness	Pay-as-you-go	Competitive pricing	Enterprise focus

📌 Choosing the Right Cloud Platform:

- ✓ Use **AWS** for flexible, scalable ML workloads.
- ✓ Use **GCP** for AI research, deep learning, and TensorFlow.
- ✓ Use **Azure** for **business-focused AI automation**.

💡 Conclusion:

AWS, GCP, and Azure each provide **unique strengths**, and the best choice depends on the **organization's data science needs**.

📌 CHAPTER 7: EXERCISES & ASSIGNMENTS

7.1 Multiple Choice Questions

- Which cloud platform is best for AI research?
- (a) AWS

- (b) GCP
- (c) Azure
- (d) IBM Cloud

Which Azure service is used for big data analytics?

- (a) Synapse Analytics
- (b) SageMaker
- (c) BigQuery
- (d) TensorFlow

7.2 Practical Assignment

 **Task:**

- Deploy a simple ML model on **AWS SageMaker, GCP AI Platform, or Azure ML.**
- Analyze big data using **BigQuery or Redshift.**
- Compare costs and performance across platforms.



DEPLOYING MACHINE LEARNING MODELS ON CLOUD

📌 CHAPTER 1: INTRODUCTION TO DEPLOYING MACHINE LEARNING MODELS

1.1 What is Model Deployment?

Machine Learning (ML) model deployment is the process of **integrating a trained model into a production environment**, where it can make real-time predictions on new data. Deployment ensures that ML models can be used by applications, businesses, or end-users to generate valuable insights.

◆ Why is Model Deployment Important?

- ✓ **Real-Time Predictions** – Once deployed, models can process and return predictions instantly.
- ✓ **Automation & Scalability** – ML models deployed on cloud platforms handle thousands to millions of requests simultaneously.
- ✓ **Integration with Applications** – ML models can power recommendation systems, fraud detection tools, chatbots, and AI assistants.
- ✓ **Data-Driven Decision Making** – Businesses can use deployed ML models to enhance customer experiences and optimize operations.

📌 Example Use Case:

An **e-commerce website** deploys an ML model on the cloud to **recommend personalized products** to customers based on their browsing history.

Conclusion:

Deploying ML models is the **final step in an AI project**, allowing real-world applications to benefit from AI-driven automation.

CHAPTER 2: CHALLENGES IN DEPLOYING MACHINE LEARNING MODELS

2.1 Why is ML Deployment Difficult?

Many ML models fail to move beyond the development phase due to deployment challenges, including:

- ✓ **Model Scalability Issues** – Handling high traffic while maintaining performance.
- ✓ **Model Performance Drift** – Over time, models become outdated due to **changing real-world data**.
- ✓ **Integration with Existing Systems** – Ensuring smooth communication with databases and APIs.
- ✓ **Security Concerns** – Protecting ML models from cyber threats and data leaks.
- ✓ **Resource Management** – Allocating **computing power** efficiently to optimize cost.

Example:

A fraud detection ML model deployed in a banking system may become **less effective over time** as fraudsters develop **new methods**. The model needs **continuous retraining** to stay relevant.

Conclusion:

Deployment is not just about launching a model—it involves **continuous monitoring, updating, and scaling** to maintain its accuracy and efficiency.

📌 CHAPTER 3: OVERVIEW OF CLOUD-BASED DEPLOYMENT

3.1 Why Deploy on Cloud Instead of Local Servers?

Cloud platforms provide **better scalability, security, and efficiency** compared to traditional on-premise servers.

◆ Benefits of Cloud Deployment:

- ✓ **High Scalability** – Handles large workloads without requiring physical hardware.
- ✓ **Automatic Model Updates** – Easily update ML models as new data arrives.
- ✓ **Efficient Resource Management** – Auto-scaling ensures **cost optimization**.
- ✓ **Easy Integration with Databases & APIs** – Connect ML models to mobile/web applications seamlessly.
- ✓ **Security & Compliance** – Cloud providers offer **built-in encryption and authentication** mechanisms.

3.2 Popular Cloud Platforms for ML Deployment

Cloud Service	Features	Example Use Cases
AWS SageMaker	Fully managed ML service	Fraud detection, personalized recommendations
Google AI Platform	Scalable ML deployment & AutoML	Chatbots, medical diagnosis
Microsoft Azure ML	Supports multiple ML frameworks	Image recognition, business forecasting

IBM Watson ML	NLP & AI-powered analytics	AI assistants, sentiment analysis
Oracle Cloud ML	High-performance computing for AI	Predictive analytics, financial forecasting

📌 **Example:**

A healthcare startup uses **Google Cloud AI Platform** to deploy a model that **analyzes medical reports** and suggests possible **diagnoses based on historical patient data**.

💡 **Conclusion:**

Choosing the right **cloud platform** depends on the **model's requirements, budget, and infrastructure compatibility**.

📌 **CHAPTER 4: STEPS TO DEPLOY ML MODELS ON CLOUD**

4.1 General Workflow for ML Model Deployment

- ✓ **Step 1: Train & Evaluate the Model** – Ensure high accuracy before deployment.
- ✓ **Step 2: Convert Model to a Deployable Format** – Save in **ONNX, TensorFlow SavedModel, or Pickle format**.
- ✓ **Step 3: Choose a Cloud Provider** – AWS, Google Cloud, or Azure.
- ✓ **Step 4: Deploy Model via API or Web Service** – Use **Flask, FastAPI, or Docker containers**.
- ✓ **Step 5: Monitor Model Performance** – Implement logging and performance tracking.

📌 **Example:**

A fintech company develops a **loan approval ML model**. They:

- ☒ Train it in **Jupyter Notebook**.
- ☒ Save it as a **.pkl file**.
- ☒ Deploy on **AWS Lambda** as an **API endpoint**.

- 4 Continuously monitor prediction accuracy to improve model performance.

Conclusion:

Deploying ML models follows a structured **workflow from training to monitoring**, ensuring smooth production integration.

CHAPTER 5: DEPLOYING ML MODELS WITH DOCKER & FASTAPI

5.1 Why Use Docker for Deployment?

Docker packages an **ML model, dependencies, and libraries** into a container that can run anywhere without compatibility issues.

- ✓ **Portability** – Runs on any machine with Docker installed.
- ✓ **Environment Consistency** – Avoids dependency conflicts.
- ✓ **Scalability** – Easily deploy multiple containers using **Kubernetes**.

5.2 Deploying an ML Model Using FastAPI & Docker

Step 1: Create a Simple FastAPI Model Endpoint

```
from fastapi import FastAPI  
import pickle  
import numpy as np  
  
app = FastAPI()
```

```
# Load ML model  
model = pickle.load(open("model.pkl", "rb"))
```

```
@app.get("/")  
def home():  
    return {"message": "Welcome to ML Model API"}
```

```
@app.post("/predict/")  
def predict(data: list):  
    prediction = model.predict(np.array(data).reshape(1, -1))  
    return {"prediction": prediction.tolist()}
```

➡ Step 2: Create a Dockerfile for Deployment

```
# Use Python base image
```

```
FROM python:3.8
```

```
# Set working directory
```

```
WORKDIR /app
```

```
# Copy required files
```

```
COPY . /app
```

```
# Install dependencies
```

```
RUN pip install -r requirements.txt
```

```
# Run FastAPI server
```

```
CMD ["uvicorn", "main:app", "--host", "0.0.0.0", "--port", "8000"]
```

📌 Step 3: Build and Run the Docker Container

```
docker build -t ml_model .
```

```
docker run -p 8000:8000 ml_model
```

💡 Conclusion:

Using **Docker** and **FastAPI** allows seamless **cloud deployment of ML models** in a lightweight, scalable format.

📌 CHAPTER 6: MONITORING DEPLOYED ML MODELS

6.1 Why Monitor Machine Learning Models?

Once deployed, an ML model needs **continuous monitoring** to ensure accuracy and efficiency.

- ✓ **Performance Drift Detection** – Detects when a model's accuracy drops.
- ✓ **Model Retraining** – Updates model with new real-world data.
- ✓ **Logging & Error Handling** – Keeps track of failed predictions.
- ✓ **Security Measures** – Prevents **unauthorized API access**.

📌 Example:

A stock price prediction model deployed on **AWS Lambda** continuously logs **prediction errors**. When accuracy drops below 85%, it **triggers automatic retraining**.

💡 Conclusion:

Deploying an ML model is **not the end**—monitoring and **continuous improvement** are required to maintain accuracy.

📌 CHAPTER 7: SUMMARY & NEXT STEPS

✓ Key Takeaways:

- ✓ Cloud platforms **AWS, GCP, and Azure** provide scalable deployment solutions.
- ✓ **FastAPI + Docker** makes cloud deployment **efficient and portable**.
- ✓ Continuous **monitoring, logging, and retraining** improve deployed ML models.

📌 Next Steps:

- ◆ Deploy an ML model using **Flask or FastAPI**.
- ◆ Host it on **AWS Lambda or Google Cloud Run**.
- ◆ Set up **model monitoring dashboards** using Grafana.

ISDM



SERVERLESS COMPUTING FOR DATA SCIENCE APPLICATIONS

📌 CHAPTER 1: INTRODUCTION TO SERVERLESS COMPUTING

1.1 What is Serverless Computing?

Serverless computing is a cloud computing execution model where the **cloud provider dynamically manages the infrastructure**, allowing developers to focus only on **writing and deploying code**. Unlike traditional cloud computing models, **serverless computing eliminates the need to manage servers, scale infrastructure, or handle provisioning**.

- ◆ **Key Characteristics of Serverless Computing:**
- ✓ **No Server Management** – Developers don't worry about provisioning or maintaining servers.
- ✓ **Scalability** – Automatically scales up or down based on demand.
- ✓ **Cost-Effective** – Users only pay for actual compute time, reducing costs.
- ✓ **Event-Driven Execution** – Functions execute **only when triggered by an event**.

📌 Example Use Case:

A data scientist deploys a **machine learning model** using AWS Lambda. The model runs **only when a user requests a prediction**, reducing operational costs compared to a traditional server setup.

💡 Conclusion:

Serverless computing **simplifies the deployment process, enhances scalability, and reduces costs**, making it ideal for **data science applications**.

📌 CHAPTER 2: KEY COMPONENTS OF SERVERLESS COMPUTING

Serverless computing consists of **various components** that help in executing applications without requiring direct server management.

2.1 Function-as-a-Service (FaaS)

- ✓ Allows developers to **run individual functions** in response to events.
- ✓ Examples: **AWS Lambda, Google Cloud Functions, Azure Functions.**

📌 Example:

An **image processing function** automatically resizes images uploaded to cloud storage using AWS Lambda.

2.2 Backend-as-a-Service (BaaS)

- ✓ Offers **ready-to-use backend services** such as authentication, databases, and API management.
- ✓ Examples: **Firebase, AWS Amplify, Supabase.**

📌 Example:

A data science application **stores user preferences and ML model outputs** in Firebase without needing a dedicated backend.

2.3 Serverless Databases

- ✓ **Fully managed databases** that scale dynamically.
- ✓ Examples: **Amazon DynamoDB, Google Firestore, Azure Cosmos DB.**

📌 Example:

A recommendation system stores **user interactions** in **DynamoDB**, which automatically scales based on user traffic.

💡 Conclusion:

Serverless computing consists of **multiple cloud services** that work together to provide a seamless, highly scalable, and cost-efficient computing environment.

📌 CHAPTER 3: SERVERLESS COMPUTING IN DATA SCIENCE

3.1 Why Use Serverless Computing in Data Science?

Serverless computing enhances **data science workflows** by automating resource management, providing scalability, and ensuring high availability.

◆ Advantages for Data Science:

- ✓ **Efficient Model Deployment** – ML models can be deployed as serverless APIs.
- ✓ **Automatic Scaling** – Increases resources dynamically when handling large datasets.
- ✓ **Cost Efficiency** – Pay only for actual compute time, ideal for batch processing.

📌 Example:

A serverless function **triggers a machine learning model** when a new dataset is uploaded to the cloud.

3.2 Common Use Cases in Data Science

◆ **Model Training & Inference:**

- ✓ ML models are trained and deployed **without dedicated servers**.
- ✓ AWS Lambda serves a **fraud detection model** on demand.

◆ **Data Preprocessing & Transformation:**

- ✓ Serverless functions clean and prepare datasets.
- ✓ Google Cloud Functions handle **ETL pipelines for big data processing**.

◆ **Real-Time Data Processing:**

- ✓ Serverless functions analyze streaming data in **IoT and financial applications**.
- ✓ Azure Functions process **real-time sensor data from industrial machines**.

💡 **Conclusion:**

Serverless computing **reduces the complexity of deploying and managing data science applications**, enabling automated and scalable workflows.

📌 CHAPTER 4: DEPLOYING MACHINE LEARNING MODELS WITH SERVERLESS COMPUTING

4.1 Steps to Deploy a Machine Learning Model Using Serverless Functions

Step 1: Train the model using **Jupyter Notebook or Google Colab**.

```
from sklearn.ensemble import RandomForestClassifier  
import pickle
```

```
# Train the model
```

```
model = RandomForestClassifier(n_estimators=100)  
model.fit(X_train, y_train)
```

```
# Save the model
```

```
pickle.dump(model, open('model.pkl', 'wb'))
```

Step 2: Create a **serverless API** using AWS Lambda and API Gateway.

Step 3: Deploy the model in a cloud storage service (**AWS S3, Google Cloud Storage**).

Step 4: Write a **serverless function** to handle incoming requests.

```
import json
```

```
import boto3
```

```
import pickle
```

```
# Load model
```

```
s3 = boto3.client('s3')
```

```
model_file = s3.get_object(Bucket='my-model-bucket',  
Key='model.pkl')
```

```
model = pickle.load(model_file['Body'])
```

```
def lambda_handler(event, context):
```

```
    input_data = json.loads(event['body'])
```

```
    prediction = model.predict([input_data['features']])
```

```
return {  
    'statusCode': 200,  
    'body': json.dumps({'prediction': prediction.tolist()})  
}
```

💡 Conclusion:

Deploying ML models with serverless computing **eliminates the need for dedicated servers** and provides **scalable, cost-efficient model serving.**

📌 CHAPTER 5: SERVERLESS DATA PIPELINES FOR DATA SCIENCE

5.1 What is a Serverless Data Pipeline?

A data pipeline **automates the process of collecting, transforming, and storing data** in a structured manner.

5.2 Components of a Serverless Data Pipeline

- ◆ **Data Ingestion (Extract)**
- ✓ AWS Lambda triggers when new data is uploaded to **S3**.
- ◆ **Data Processing (Transform)**
- ✓ Google Cloud Functions process data into structured formats.
- ◆ **Data Storage (Load)**
- ✓ Data is stored in **BigQuery or DynamoDB** for analysis.

📌 Example:

A serverless pipeline automatically **ingests social media data, cleans it, and stores it in a data warehouse** for real-time analysis.

Conclusion:

Serverless data pipelines **simplify data engineering tasks**, making real-time data processing efficient and scalable.

CHAPTER 6: ADVANTAGES AND CHALLENGES OF SERVERLESS COMPUTING

6.1 Advantages of Serverless Computing for Data Science

- ✓ **Scalability:** Serverless functions **automatically scale** based on workload.
- ✓ **Cost Savings:** You **only pay for execution time**, making it more cost-effective.
- ✓ **Faster Deployment:** No need to set up and configure servers.
- ✓ **Improved Fault Tolerance:** Cloud providers handle **redundancy and fault recovery**.

6.2 Challenges and Limitations

- ✓ **Cold Start Latency:** Initial execution may have a delay.
- ✓ **Limited Execution Time:** Most serverless functions have **timeout** limits.
- ✓ **Not Ideal for Long Running Tasks:** Best suited for **short-lived tasks** rather than continuous processes.

Conclusion:

Despite challenges, serverless computing **offers significant benefits** for data science applications, especially in terms of **cost efficiency, scalability, and automation**.

📌 CHAPTER 7: FUTURE OF SERVERLESS COMPUTING IN DATA SCIENCE

7.1 Trends in Serverless Computing

- ✓ **Integration with AI & ML:** Serverless platforms are now supporting **GPU-powered ML workloads**.
- ✓ **Serverless Containers:** New services like **AWS Fargate** allow running serverless containerized workloads.
- ✓ **Edge Computing & IoT:** Serverless is being used for **real-time IoT data processing**.

📌 Example:

Tesla's **autonomous driving data pipeline** processes car sensor data using **serverless cloud functions** for **real-time decision-making**.

💡 Conclusion:

The future of serverless computing is **deeply intertwined with AI and big data**, making it an essential technology for **scalable, real-time, and cost-efficient data science workflows**.

📌 FINAL THOUGHTS & NEXT STEPS

✓ Key Takeaways:

- ✓ Serverless computing **eliminates the need for managing servers**.
- ✓ Data science applications **benefit from auto-scaling, cost efficiency, and rapid deployment**.
- ✓ Serverless **ML model deployment** simplifies inference at scale.
- ✓ Serverless **data pipelines** automate ETL workflows.

📌 Next Steps:

- ◆ Experiment with **AWS Lambda, Google Cloud Functions, and Azure Functions**.
- ◆ Deploy a **real-world machine learning model using a serverless**

API.

- ◆ Build a **serverless ETL pipeline for streaming data**.

ISDM-NxT



APACHE KAFKA FOR REAL-TIME DATA PROCESSING

📌 CHAPTER 1: INTRODUCTION TO APACHE KAFKA

1.1 What is Apache Kafka?

Apache Kafka is a **distributed event streaming platform** designed for **high-throughput, real-time data processing**. It allows applications to **publish, subscribe, store, and process data streams in real-time** across multiple systems efficiently.

Kafka was originally developed at **LinkedIn** and later open-sourced under **Apache Software Foundation**. Today, it is widely used by **large-scale organizations** like Netflix, Uber, Twitter, and LinkedIn for **real-time analytics, monitoring, event-driven architectures, and log aggregation**.

◆ Why Use Apache Kafka?

- ✓ **High Scalability** – Handles millions of messages per second.
- ✓ **Fault Tolerant & Reliable** – Ensures no data loss, even if brokers fail.
- ✓ **Real-time Streaming** – Processes data as events occur.
- ✓ **Decouples Data Producers & Consumers** – Enhances system flexibility.
- ✓ **Distributed & Resilient** – Spreads load across multiple nodes.

📌 Example Use Case:

An **e-commerce platform** uses Kafka to **track customer actions (clicks, purchases, searches) in real-time** and deliver personalized recommendations instantly.

Conclusion:

Apache Kafka is a **powerful event-driven architecture** that provides **high-speed, fault-tolerant, and scalable** messaging solutions for modern applications.

CHAPTER 2: ARCHITECTURE OF APACHE KAFKA

Kafka follows a **distributed, publish-subscribe messaging system** where data moves as real-time streams.

2.1 Core Components of Kafka

- ✓ **Producer** – Sends messages to Kafka topics.
- ✓ **Broker** – Stores and distributes messages to consumers.
- ✓ **Consumer** – Reads messages from Kafka topics.
- ✓ **Topic** – A category where messages are published.
- ✓ **Partition** – A way to split a topic for parallel processing.
- ✓ **ZooKeeper** – Manages Kafka clusters and leader election.

Example:

A **weather monitoring system** collects temperature data from sensors (Producers) and streams it to multiple weather forecasting applications (Consumers) through Kafka topics.

2.2 How Kafka Works?

- 1 Producers send messages to Kafka topics.
- 2 Messages are stored in partitions across multiple Kafka brokers.
- 3 Consumers subscribe to topics and read messages in order.
- 4 Zookeeper maintains metadata and leader-follower relationships for partitions.

💡 Conclusion:

Kafka's **distributed architecture** ensures **high availability and efficient event streaming**, making it ideal for real-time processing.

📌 CHAPTER 3: INSTALLING AND CONFIGURING APACHE KAFKA

3.1 Setting Up Kafka Locally (Linux/Windows/Mac)

📌 Install Kafka Using Binary Downloads:

Download Kafka

```
 wget https://downloads.apache.org/kafka/3.0.0/kafka_2.13-3.0.0.tgz
```

```
 tar -xvzf kafka_2.13-3.0.0.tgz
```

```
 cd kafka_2.13-3.0.0
```

📌 Start Zookeeper Service:

```
 bin/zookeeper-server-start.sh config/zookeeper.properties
```

📌 Start Kafka Server:

```
 bin/kafka-server-start.sh config/server.properties
```

📌 Create a Kafka Topic:

```
 bin/kafka-topics.sh --create --topic test_topic --bootstrap-server localhost:9092 --partitions 3 --replication-factor 1
```

📌 Send & Receive Messages:

Start producer

```
 bin/kafka-console-producer.sh --topic test_topic --bootstrap-server localhost:9092
```

> Hello Kafka

```
# Start consumer
```

```
bin/kafka-console-consumer.sh --topic test_topic --from-beginning -  
-bootstrap-server localhost:9092
```

 **Conclusion:**

Kafka is easy to set up on **local machines or cloud environments** for development and testing.

 **CHAPTER 4: PRODUCER & CONSUMER APIs IN KAFKA**

Kafka provides **high-performance APIs** for integrating applications.

4.1 Kafka Producer API (Sending Messages)

Kafka Producers **push data to Kafka topics asynchronously**.

 **Example (Python Producer Using Kafka-Python Library):**

```
from kafka import KafkaProducer
```

```
producer = KafkaProducer(bootstrap_servers='localhost:9092')
```

```
producer.send('test_topic', b'Hello, Kafka!')
```

```
producer.flush()
```

 **Example (Java Producer):**

```
Properties props = new Properties();
```

```
props.put("bootstrap.servers", "localhost:9092");
```

```
props.put("key.serializer",
```

```
"org.apache.kafka.common.serialization.StringSerializer");
```

```
props.put("value.serializer",
"org.apache.kafka.common.serialization.StringSerializer");

KafkaProducer<String, String> producer = new
KafkaProducer<>(props);

producer.send(new ProducerRecord<>("test_topic", "Hello Kafka!"));

producer.close();
```

4.2 Kafka Consumer API (Reading Messages)

Kafka Consumers **subscribe to topics** and **consume messages in real-time**.

📌 **Example (Python Consumer):**

```
from kafka import KafkaConsumer

consumer = KafkaConsumer('test_topic',
bootstrap_servers='localhost:9092')

for message in consumer:
    print(message.value.decode())
```

📌 **Example (Java Consumer):**

```
KafkaConsumer<String, String> consumer = new
KafkaConsumer<>(props);

consumer.subscribe(Collections.singletonList("test_topic"));

while (true) {
```

```
ConsumerRecords<String, String> records =  
consumer.poll(Duration.ofMillis(100));  
  
for (ConsumerRecord<String, String> record : records) {  
  
    System.out.println("Received: " + record.value());  
  
}  
  
}
```

💡 Conclusion:

Kafka's Producer and Consumer APIs **enable seamless real-time data exchange** across systems.

📌 CHAPTER 5: KAFKA STREAM PROCESSING WITH KAFKA STREAMS & FLINK

Kafka enables **real-time stream processing** using tools like **Kafka Streams, Apache Flink, and Spark Streaming**.

5.1 Kafka Streams (Processing Kafka Data in Real-time)

Kafka Streams is a **Java-based framework** for processing streaming data.

📌 Example (Kafka Stream Processing in Java):

```
StreamsBuilder builder = new StreamsBuilder();  
  
KStream<String, String> source = builder.stream("input_topic");  
  
KStream<String, String> processed = source.mapValues(value ->  
value.toUpperCase());  
  
processed.to("output_topic");  
  
KafkaStreams streams = new KafkaStreams(builder.build(),  
properties);
```

```
streams.start();
```

💡 Conclusion:

Kafka Streams allows **real-time transformation of streaming data** before it reaches consumers.

📌 CHAPTER 6: MONITORING, PERFORMANCE TUNING & SCALABILITY IN KAFKA

6.1 Monitoring Kafka with Prometheus & Grafana

✓ Kafka exposes metrics via JMX that can be monitored using Prometheus & Grafana dashboards.

📌 Example (Enable Kafka Metrics):

```
export KAFKA_OPTS="-javaagent:/path/to/jmx_exporter.jar"
```

✓ Use Prometheus for collecting Kafka logs and Grafana for visualizing broker performance.

6.2 Performance Tuning Kafka for High Throughput

✓ Increase Producer Throughput:

```
batch.size=16384
```

```
linger.ms=5
```

```
compression.type=lz4
```

✓ Optimize Consumer Lag Reduction:

```
fetch.min.bytes=50000
```

```
fetch.max.wait.ms=100
```

✓ Scaling Kafka Brokers for High Availability:

- Deploy Kafka in **multi-node clusters** across different regions.
- Use **horizontal scaling** by adding more partitions.

💡 Conclusion:

Proper **monitoring and tuning** ensures Kafka can handle **millions of real-time events per second** without bottlenecks.

📌 CHAPTER 7: SUMMARY & NEXT STEPS

✓ Key Takeaways:

- ✓ Kafka is a **highly scalable** real-time data streaming platform.
- ✓ Producers publish messages, while consumers subscribe to topics.
- ✓ Kafka Streams and Apache Flink enable **real-time data processing**.
- ✓ Cloud-native monitoring solutions like **Prometheus and Grafana** optimize performance.

📌 Next Steps:

- ◆ Deploy Kafka in a cloud environment (AWS, Azure, Google Cloud).
- ◆ Build a real-time analytics system using Kafka & Apache Flink.
- ◆ Explore Kafka integration with Spark Streaming & Elasticsearch.

📌 **ASSIGNMENT:**
☑ **DEPLOY A MACHINE LEARNING MODEL
AS A REST API USING FLASK & AWS**

ISDM-Nxt

🔧 SOLUTION: DEPLOYING A MACHINE LEARNING MODEL AS A REST API USING FLASK & AWS

🎯 Objective:

The goal of this assignment is to **train a Machine Learning model, deploy it as a REST API using Flask, and host it on AWS**. This allows users to send input data via an API request and receive predictions from the model.

◆ STEP 1: Train and Save the Machine Learning Model

We first train a **Machine Learning model (e.g., a classification model)** and save it for deployment.

1.1 Install Required Libraries

```
pip install pandas numpy scikit-learn flask gunicorn boto3
```

1.2 Train and Save the Model (Logistic Regression Example)

```
import pandas as pd  
import numpy as np  
from sklearn.model_selection import train_test_split  
from sklearn.linear_model import LogisticRegression  
from sklearn.metrics import accuracy_score  
import pickle
```

```
# Load dataset
```

```
df = pd.read_csv("customer_churn.csv")

# Selecting features and target variable
X = df[['tenure', 'monthly_charges', 'total_charges']]
y = df['churn']

# Splitting into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Train a logistic regression model
model = LogisticRegression()
model.fit(X_train, y_train)

# Evaluate model
y_pred = model.predict(X_test)
print(f"Model Accuracy: {accuracy_score(y_test, y_pred):.2f}")

# Save the trained model as a pickle file
with open("model.pkl", "wb") as f:
    pickle.dump(model, f)
```

🔍 Key Insights:

- **Pickle** saves the model so it can be loaded into a Flask API without retraining.
 - The model predicts **customer churn** based on input features like tenure and monthly charges.
-

◆ STEP 2: Build a Flask REST API

Flask is a **lightweight web framework** for building RESTful APIs.

2.1 Create a Flask Application

Create a new file **app.py**:

```
from flask import Flask, request, jsonify  
import pickle  
import numpy as np
```

```
# Load trained model  
with open("model.pkl", "rb") as f:  
    model = pickle.load(f)  
  
# Initialize Flask app  
app = Flask(__name__)
```

```
# Define a route for prediction  
@app.route('/predict', methods=['POST'])  
def predict():
```

```
data = request.get_json() # Get JSON input from user  
  
features = np.array([data['tenure'], data['monthly_charges'],  
data['total_charges']]).reshape(1, -1)  
  
prediction = model.predict(features)[0] # Make prediction  
  
return jsonify({"Churn Prediction": "Yes" if prediction == 1 else  
"No"})
```

```
# Run Flask app  
  
if __name__ == '__main__':  
  
    app.run(debug=True, host='0.0.0.0', port=5000)
```

2.2 Explanation of Code

- ✓ `request.get_json()` – Accepts JSON input from API users.
 - ✓ `np.array()` – Converts JSON data into a format suitable for model prediction.
 - ✓ `model.predict()` – Uses the saved model to make predictions.
 - ✓ `jsonify()` – Returns the prediction result in JSON format.
-

◆ STEP 3: Test the Flask API Locally

3.1 Run the Flask App

```
python app.py
```

It will start a local web server at <http://127.0.0.1:5000/>.

3.2 Test API Using Postman or Curl

Using Postman:

1. Open Postman and select **POST** request.

2. Enter URL: `http://127.0.0.1:5000/predict`

3. In **Body → raw → JSON**, send:

```
{  
    "tenure": 12,  
    "monthly_charges": 50.0,  
    "total_charges": 600.0  
}
```

4. Click **Send** to get the response:

```
{  
    "Churn Prediction": "No"  
}
```

Using Curl Command:

```
curl -X POST http://127.0.0.1:5000/predict -H "Content-Type: application/json" -d '{"tenure": 12, "monthly_charges": 50.0, "total_charges": 600.0}'
```

🔍 Key Insights:

- The API correctly processes input data and returns a churn prediction.
- Debugging locally ensures the API works correctly before deployment.

◆ STEP 4: Deploy Flask App on AWS

4.1 Setup an AWS EC2 Instance

-
- 1 Sign in to AWS Console → Navigate to EC2.
 - 2 Click Launch Instance → Choose Ubuntu 20.04.
 - 3 Choose t2.micro (Free Tier) → Click Next.
 - 4 Configure storage (default 8GB) → Click Next.
 - 5 Add security group → Allow HTTP, HTTPS, and port 5000.
 - 6 Click Launch → Download key pair (.pem) file.
-

4.2 Connect to EC2 via SSH

Use the **key pair** to connect to the EC2 instance.

```
ssh -i "your-key.pem" ubuntu@your-ec2-public-ip
```

4.3 Install Python & Dependencies

```
sudo apt update
```

```
sudo apt install python3-pip python3-venv -y
```

```
pip3 install flask scikit-learn numpy pandas gunicorn
```

4.4 Upload Flask App & Model to EC2

Use SCP to upload files from your local machine to AWS EC2.

```
scp -i "your-key.pem" app.py model.pkl ubuntu@your-ec2-public-ip:~/flask-app/
```

4.5 Run Flask App on EC2

```
cd flask-app
```

```
python3 app.py
```

- Your API will be available at <http://your-ec2-public-ip:5000/predict>.
 - Use **Postman** or **Curl** to test the API like before.
-

◆ STEP 5: Deploy API using Gunicorn & Nginx

5.1 Install Gunicorn & Run Flask App in Background

```
pip3 install gunicorn
```

```
gunicorn --bind 0.0.0.0:5000 app:app
```

5.2 Install & Configure Nginx for Reverse Proxy

```
sudo apt install nginx -y
```

```
sudo nano /etc/nginx/sites-available/flask_app
```

Add the following configuration:

```
server {  
    listen 80;  
    server_name your-ec2-public-ip;  
  
    location / {  
        proxy_pass http://127.0.0.1:5000;  
        proxy_set_header Host $host;  
        proxy_set_header X-Real-IP $remote_addr;  
    }  
}
```

Save and enable the configuration:

```
sudo ln -s /etc/nginx/sites-available/flask_app /etc/nginx/sites-enabled/
```

```
sudo systemctl restart nginx
```

Your API is now available at <http://your-ec2-public-ip/>.

➡ STEP 6: Testing the Live API on AWS

```
curl -X POST http://your-ec2-public-ip/predict -H "Content-Type: application/json" -d '{"tenure": 12, "monthly_charges": 50.0, "total_charges": 600.0}'
```

Expected Output:

```
{  
    "Churn Prediction": "No"  
}
```

➡ SUMMARY OF STEPS

- Trained a machine learning model for customer churn prediction.
 - Created a Flask API to serve model predictions.
 - Tested API locally using Postman and Curl.
 - Deployed API on AWS EC2 instance with Gunicorn & Nginx.
-

📌 NEXT STEPS

📌 Enhancements:

- ◆ Store customer request logs in a **database** for monitoring.
- ◆ Implement **authentication & security** in API using JWT tokens.
- ◆ Containerize the API using **Docker** for portability.

ISDM-Nxt