# CS 643 Cloud Computing: Programming Assignment 2

## Energy Consumption Forecasting with Apache Spark and Docker

Hemanth Sruthi Vellore [HV234]

May 6, 2025

**Abstract:** This report presents an energy consumption forecasting system developed for CS 643 Cloud Computing Programming Assignment 2, leveraging Apache Spark on Amazon AWS. The project encompasses parallel model training across multiple EC2 instances, a streamlined single-instance prediction application, and a Docker container for seamless deployment. Powered by a Gradient Boosted Trees regression model from Spark MLlib, the system is trained on `TrainingDataset.csv` and validated with `ValidationDataset.csv`. This document provides a comprehensive guide for configuring the cloud infrastructure, executing the training and prediction workflows, and deploying the Docker container. It also evaluates the role of AI tools (ChatGPT) in code development, offering attribution and insights into the development process.

# Contents

# 1   Technical Architecture

The system is structured around three key components: distributed model training, single-instance prediction, and Docker-based deployment.
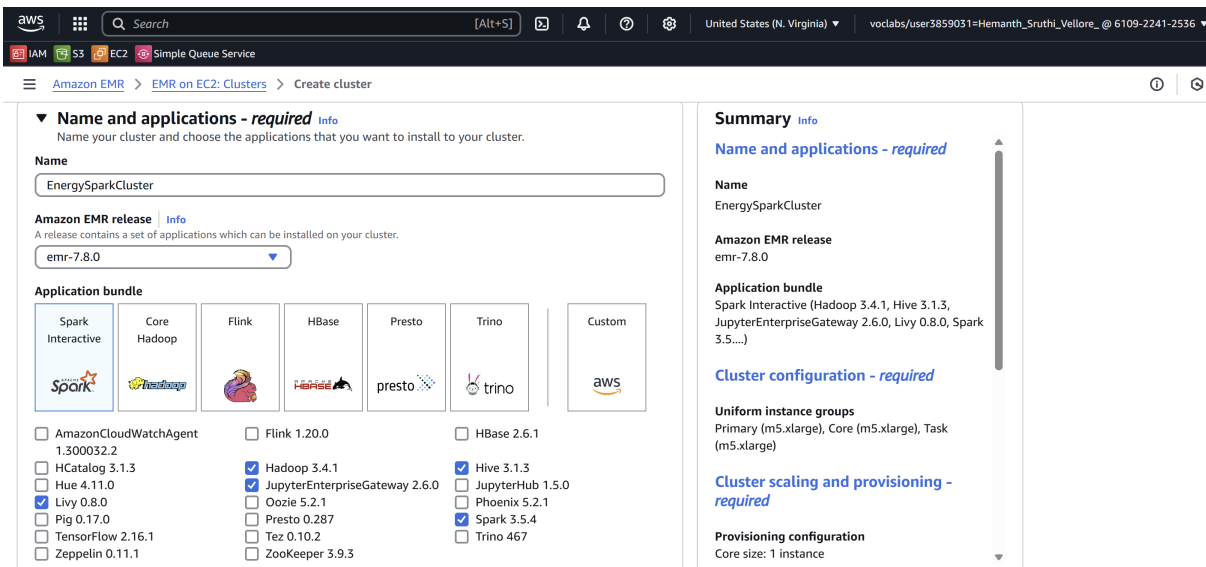
## 1.1   Distributed Model Training

This component utilizes Apache Spark on a cluster of four EC2 instances to train a GBT regression model on `TrainingDataset.csv`.

### 1.1.1   Cluster Provisioning

A Spark cluster is provisioned on AWS EC2 instances using the following steps:

1. From the AWS console, go to EMR Service and click Create Cluster.

2. Configure the cluster as follows:

   (a) Enter cluster name EnergySparkCluster

   (b) Select Release emr-7.8.0

   (c) Choose Spark Interactive: Spark 3.5.4, Hadoop 3.4.1



   (d) Configure Hardware Configurations:

       i. Select OS option (Amazon Linux release).

      ii. Select the instance type (m5.xlarge).

iii. Set number of instances to 4 (1 master, 3 slaves).

(e) Select the "vockey " key pair

(f) Select bootstrap script from the S3 bucket.



3. Click on Create Cluster.



The cluster is initialized with a bootstrap script to install dependencies.

**Listing 1:** Bootstrap Script (`boot.sh`)

```
sudo pip3 install numpy pandas
```

## 1.1.2 Uploading Files to the EMR Cluster Master Node

After the EMR cluster has been successfully provisioned, transfer the required files to the master node using a secure file transfer client such as WinSCP. Establish a connection to the master node with the appropriate SSH credentials. Upload the `TrainingDataset.csv` file and the `energy_model_trainer.py` script to the home directory of the `hadoop` user.

### 1.1.3   HDFS Data Integration

Now that all the files are on the master node, we want to migrate them to HDFS so that all the slave nodes can access them without having to physically copy them to all EC2 nodes. SSH into the master node and run the below commands to copy the files to HDFS.

**Listing 2:** HDFS Integration Commands

```
hadoop fs -put TrainingDataset.csv /user/hadoop/TrainingDataset.csv
hadoop fs -put energy_model_trainer.py /user/hadoop/energy_model_trainer.py
hdfs dfs -ls -t -R
```

**Figure 1:** Data integration and verification.



### 1.1.4   Model Training Workflow

The `energy_model_trainer.py` script, executed via `spark-submit`, trains the GBT model by loading the dataset, encoding categorical features, assembling feature vectors, splitting data into training and testing sets, training the model, evaluating performance (RMSE and $R^2$), and saving the model to `EnergyPredictorGBT`.

**Listing 3:** Training Workflow Command

```
spark-submit energy_model_trainer.py
```

### 1.1.5 Training Performance Metrics

The model achieved a Root Mean Squared Error (RMSE) of 70.4804 and a coefficient of determination ($R^2$) of 0.9942 on the training dataset, indicating a strong fit and high predictive accuracy.



### 1.1.6 Model Export Process

The trained model was exported to the master node's local file system, compressed into an archive, and downloaded to the local machine using WinSCP. After verifying the download, the EMR cluster was shut down to avoid unnecessary resource usage.

**Listing 4:** Model Export Commands

```
hdfs dfs -copyToLocal EnergyPredictorGBT /home/hadoop/EnergyPredictorGBT
tar -czf model.tar.gz EnergyPredictorGBT/
```

## 1.2    Single-Instance Prediction System

The prediction system runs on a single EC2 instance using `energy_usage_predictor.py` and the saved model.

### 1.2.1    EC2 Instance Creation

Provisioned an EC2 instance via the AWS Management Console by:

- Naming the instance `EnergyPredictor-EC2`

- Selecting `Ubuntu Server 24.04 LTS (HVM)` AMI

- Choosing `t2.medium` instance type

- Setting root volume to 10 GB (gp3)

- Assigning the `EMR_EC2_DefaultRole` IAM role

- Selecting the `vockey` key pair

- Launching the instance

Instance initialized and ready for setup.

### 1.2.2  EC2 Instance Setup

After creating the instance, SSH access is established using the key pair. The system is updated, and essential dependencies, including Java, Python, and Apache Spark, are installed and configured to prepare the environment.

**Listing 5:** Instance Setup Commands

```
1  sudo apt-get update
2  sudo apt-get install -y python3-pip python3-numpy python3-pandas openjdk-11-
     jdk
3  wget https://archive.apache.org/dist/spark/spark-3.5.5/spark-3.5.5-bin-
     hadoop3.tgz
4  sudo tar xvf spark-3.5.5-bin-hadoop3.tgz -C /opt
5  sudo chown -R ubuntu:ubuntu /opt/spark-3.5.5-bin-hadoop3
6  sudo ln -fs spark-3.5.5-bin-hadoop3 /opt/spark
```

### 1.2.3  Environment Configuration

Environment variables are set for Spark and Java.

**Listing 6:** Environment Configuration (`~/.predictor_env`)

```
1  export SPARK_HOME=/opt/spark
2  PATH=$PATH:$SPARK_HOME/bin
3  export PATH
4  export JAVA_HOME=/usr/lib/jvm/java-1.11.0-openjdk-amd64
5  export PATH=$JAVA_HOME/bin:$PATH
6  source ~/.predictor_env
```

## 1.2.4 Spark Logging Configuration

To reduce console output during execution, adjust the default Spark logging level:

- Copy the template logging configuration file.

- Edit the `log4j2.properties` file, changing `rootLogger.level` from `INFO` to `ERROR`.

**Listing 7:** Logging Configuration

```
1  cp $SPARK_HOME/conf/log4j2.properties.template $SPARK_HOME/conf/log4j2.
       properties
2  nano $SPARK_HOME/conf/log4j2.properties
3  # Change rootLogger.level = info to ERROR
```

## 1.2.5 Load Model and Script

Using WinSCP, transfer the previously exported model and the `energy_usage_predictor.py` script to the EC2 instance.

### 1.2.6 Unzip the Model

Once the files are transferred, unzip the model archive using the following command:

**Listing 8:** Unzipping the Model

```
tar -xzvf model.tar.gz
```



```
ubuntu@ip-172-31-26-90: ~
ubuntu@ip-172-31-26-90:~$ tar -xzvf model.tar.gz
EnergyPredictorGBT/
EnergyPredictorGBT/data/
EnergyPredictorGBT/data/_SUCCESS
EnergyPredictorGBT/data/part-00000-62067a7c-570b-474c-b30d-4487c8536d5b-c000.snappy.parquet
EnergyPredictorGBT/metadata/
EnergyPredictorGBT/metadata/_SUCCESS
EnergyPredictorGBT/metadata/part-00000
```

### 1.2.7 Run the Prediction Script

Execute the `energy_usage_predictor.py` script using `spark-submit` to generate predictions from the validation dataset. The script loads the dataset, preprocesses it, loads the GBT model, generates predictions, and reports RMSE and $R^2$ metrics. The results for the validation set are:

- RMSE: 73.8957

- $R^2$: 0.9933

**Listing 9:** Prediction Execution Commands

```
spark-submit energy_usage_predictor.py ValidationDataset.csv
```



```
ubuntu@ip-172-31-26-90: ~
ubuntu@ip-172-31-26-90:~$ ls
EnergyPredictorGBT  ValidationDataset.csv  energy_usage_predictor.py  model.tar.gz  spark-3.5.5-bin-hadoop3.tgz
ubuntu@ip-172-31-26-90:~$ spark-submit energy_usage_predictor.py ValidationDataset.csv
================================================================================
                        ENERGY CONSUMPTION PREDICTOR
================================================================================
A Spark-powered application for estimating energy usage on unseen datasets,
leveraging a pre-trained Gradient Boosted Trees regression model.

[STEP] Setting up Spark environment........................................... Completed
[STEP] Importing test data.................................................... Completed
[STEP] Standardizing column names............................................. Completed
[STEP] Transforming categorical features...................................... Completed
[STEP] Removing original categorical fields................................... Completed
[STEP] Casting columns to float............................................... Completed
[STEP] Preparing features and target variable................................. Completed
[STEP] Converting DataFrame to RDD format..................................... Completed
[STEP] Retrieving pre-trained model........................................... Completed
[STEP] Generating predictions................................................. Completed
[STEP] Aligning predictions with true values.................................. Completed
[STEP] Evaluating prediction performance...................................... Completed


================================================================================
                        Prediction Performance Report
================================================================================
Root Mean Squared Error (RMSE)                    73.8957
R2 (Coefficient of Determination)                  0.9933
================================================================================


================================================================================
                   ENERGY CONSUMPTION PREDICTOR COMPLETED
================================================================================
[STEP] Terminating Spark environment.......................................... Completed
ubuntu@ip-172-31-26-90:~$
```

## 1.3 Docker-Enabled Deployment

A Docker container is built to facilitate portable deployment of the prediction system.

### 1.3.1 Dockerfile Configuration

The `Dockerfile` sets up the environment and installs dependencies to run the Spark application.

- Base Image: Uses `ubuntu:20.04`.

- Dependencies: Installs Python 3, Java 11, Spark, and necessary tools.

- Environment Variables: Configures `JAVA_HOME` and `SPARK_HOME`.

- Copy Files: Transfers `energy_usage_predictor.py` and model files to the container.

- Entry Point: Uses `spark-submit` to run the script.

**Listing 10:** Dockerfile

```
1  FROM ubuntu:20.04
2  ENV DEBIAN_FRONTEND=noninteractive
3  WORKDIR /app
4  RUN apt-get update && \
5      apt-get install -y python3 python3-pip python3-numpy python3-pandas
           openjdk-11-jdk wget nano && \
6      rm -rf /var/lib/apt/lists/*
7  ENV JAVA_HOME=/usr/lib/jvm/java-1.11.0-openjdk-amd64
8  ENV PATH=$JAVA_HOME/bin:$PATH
9  RUN wget https://archive.apache.org/dist/spark/spark-3.5.5/spark-3.5.5-bin-
       hadoop3.tgz && \
10     tar xvf spark-3.5.5-bin-hadoop3.tgz -C /opt && \
11     rm spark-3.5.5-bin-hadoop3.tgz && \
12     ln -fs /opt/spark-3.5.5-bin-hadoop3 /opt/spark
13 ENV SPARK_HOME=/opt/spark
14 ENV PATH=$PATH:$SPARK_HOME/bin
15 RUN cp $SPARK_HOME/conf/log4j2.properties.template $SPARK_HOME/conf/log4j2.
       properties && \
16     sed -i 's/rootLogger.level = info/rootLogger.level = ERROR/g'
           $SPARK_HOME/conf/log4j2.properties
17 COPY energy_usage_predictor.py /app/
18 COPY EnergyPredictorGBT /app/EnergyPredictorGBT
19 ENTRYPOINT ["spark-submit", "energy_usage_predictor.py"]
20 CMD [""]
```

### 1.3.2 Container Build and Deployment

The Docker image is built, executed, and deployed as follows:

- Install Docker: Install Docker on the EC2 instance.

**Listing 11:** Install Docker

```
1    sudo apt-get install docker.io
```

- Build Image: Build the Docker image with the tag `energy-predict`.

**Listing 12:** Build Docker Image

```
1    sudo docker build -t energy-predict .
```

- Run Container: Run the container with the validation dataset mounted.

**Listing 13:** Run Docker Container

```
1    sudo docker run -v /home/ubuntu/ValidationDataset.csv:/app/
        ValidationDataset.csv energy-predict /app/ValidationDataset.csv
```

```
ubuntu@ip-172-31-26-90:~$ sudo docker run -v /home/ubuntu/ValidationDataset.csv:/app/ValidationDataset.csv energy-predict /app/ValidationDataset.csv
================================================================
                    ENERGY CONSUMPTION PREDICTOR
================================================================
A Spark-powered application for estimating energy usage on unseen datasets,
leveraging a pre-trained Gradient Boosted Trees regression model.

[STEP] Setting up Spark environment........................................... Completed
[STEP] Importing test data.................................................... Completed
[STEP] Standardizing column names............................................. Completed
[STEP] Transforming categorical features...................................... Completed
[STEP] Removing original categorical fields................................... Completed
[STEP] Casting columns to float............................................... Completed
[STEP] Preparing features and target variable................................. Completed
[STEP] Converting DataFrame to RDD format..................................... Completed
[STEP] Retrieving pre-trained model........................................... Completed
[STEP] Generating predictions................................................. Completed
[STEP] Aligning predictions with true values.................................. Completed
[STEP] Evaluating prediction performance...................................... Completed


================================================================
                    Prediction Performance Report
================================================================
Root Mean Squared Error (RMSE)                    73.8957
R2 (Coefficient of Determination)                  0.9933
================================================================


================================================================
                ENERGY CONSUMPTION PREDICTOR COMPLETED
================================================================
[STEP] Terminating Spark environment.......................................... Completed
ubuntu@ip-172-31-26-90:~$
```

- Login, Tag, and Push: Login to Docker Hub, tag the image, and push it to Docker Hub.

**Listing 14:** Login, Tag, and Push Docker Image

```
1    sudo docker login
2    sudo docker tag energy-predict sruthivellore/energy-predict
3    sudo docker push sruthivellore/energy-predict
```

- Pull Image: Pull the image from Docker Hub for use elsewhere.

**Listing 15:** Pull Docker Image

```
sudo docker pull sruthivellore/energy-predict
```

# 2   Deployment Guide

This section provides a comprehensive guide to replicating the project, ensuring clarity and reproducibility.

## 2.1   Setting Up Distributed Training

1. Provision four EC2 instances on AWS using tools (EMR).

2. Apply the bootstrap script (`boot.sh`) during instance setup.

3. Upload `TrainingDataset.csv` and `energy_model_trainer.py` to HDFS.

4. Execute `spark-submit energy_model_trainer.py`.

5. Export the trained model using `hdfs dfs -copyToLocal` and `tar`.

## 2.2   Configuring the Prediction System

1. Launch a single EC2 instance running Ubuntu 20.04.

2. Install dependencies and Spark

3. Configure environment variables (`~/.predictor_env`) and logging settings.

4. Transfer `model.tar.gz`, `energy_usage_predictor.py`, and `ValidationDataset.csv`.

5. Extract the model and run `spark-submit energy_usage_predictor.py ValidationDataset.csv`.

## 2.3   Deploying with Docker

1. Install Docker on an EC2 instance.

2. Build the Docker image using the provided `Dockerfile`.

3. Run the container, mapping the validation dataset.

4. Optionally, pull the image from Docker Hub: `sruthivellore/energy-predict`.

# 3 AI-Assisted Development

For this project, I used ChatGPT mainly to get some basic templates for scripts and setup steps. It helped with getting started faster, especially for things like the Dockerfile and a few bash commands. However, most of the actual code, including model training, prediction scripts, environment setup, and AWS configuration, was done manually based on what the project needed.

While ChatGPT saved time on some repetitive parts, a lot of debugging, fixing compatibility issues with Spark, and making everything work properly had to be done manually. Overall, AI tools were helpful for speeding up small parts.

# 4   Code and Resource Links

- GitHub Repository: `https://github.com/sruthivellore/EnergyPredictSparkAWS.git`

- Docker Hub: `https://hub.docker.com/r/sruthivellore/energy-predict`